

Improving the Numerical Performance of BLP Static and Dynamic Discrete Choice Random Coefficients Demand Estimation*

Jean-Pierre Dubé

Graduate School of Business

University of Chicago

Jeremy T. Fox

Department of Economics

University of Chicago and

NBER

Che-Lin Su

Graduate School of Business

University of Chicago

October 2008

Abstract

The widely-used estimator of Berry, Levinsohn and Pakes (1995) produces consistent instrumental variables estimates of consumer preferences from a discrete-choice demand model with random coefficients, market-level demand shocks and potentially endogenous regressors (prices). The nested fixed-point algorithm typically used for estimation is computationally intensive, largely because a system of market share equations must be repeatedly numerically inverted. We provide numerical theory results that characterize the properties of typical nested fixed-point implementations. We use these results to discuss several problems with typical computational implementations and, in particular, cases which can lead to incorrect parameter estimates. As a solution, we introduce a new computational formulation of the estimator that recasts estimation as a mathematical program with equilibrium constraints (MPEC). In many instances, MPEC is faster than the nested fixed point approach. It also avoids the numerical issues associated with nested inner loops. Several Monte Carlo experiments support our numerical concerns about NFP and the advantages of MPEC. We also discuss estimating static BLP using maximum likelihood instead of GMM. Finally, we show that MPEC is particularly attractive for forward-looking demand models where both Bellman's equation and the market share equations must be repeatedly solved.

*We thank John Birge, Lars Hansen, Kyoo il Kim, Kenneth Judd, Sven Leyffer, Aviv Nevo, Jorge Nocedal, Hugo Salgado and Richard Waltz for helpful discussions and comments. Thanks to workshop participants at Chicago, INFORMS, the International Industrial Organization Conference, Northwestern, the Portuguese Competition Commission, Rochester and the Stanford Institute for Theoretical Economics. Dubé is grateful to the Kilts Center for Marketing and the Neubauer Faculty Fund for research support. Fox thanks the NSF, grant 0721036, the Olin Foundation, and the Stigler Center for financial support. Su is grateful for the financial support from the NSF (award no. SES-0631622) and the Chicago GSB. Our email addresses are jdube@chicagogsb.edu, fox@uchicago.edu and Che-Lin.Su@chicagogsb.edu.

1 Introduction

The discrete choice class of demand models has become popular in the demand estimation literature due to the models' ability to accommodate rich substitution patterns between a potentially large array of products. The simulated method of moments estimator developed in Berry, Levinsohn and Pakes (1995), hereafter BLP, made an important contribution to this literature by accommodating controls for the endogeneity of product characteristics (namely prices) without sacrificing the flexibility of these substitution patterns. BLP consider a random coefficients discrete choice model with market-level demand shocks that correlate with prices. They construct moment conditions with which they can address the price endogeneity using standard instrumental variables methods. The approach has had a large impact: as of October 2008, BLP generated of 1000 citations in Google Scholar and the approach has been used in many important empirical studies. However, the estimator is difficult to program and can take a long time to run on a desktop computer. More importantly, some current implementations of the estimator are sufficiently vulnerable to numerical inaccuracy that they may produce incorrect parameter estimates. We summarize some of these computational problems and propose an alternative procedure that is robust to these sources of numerical inaccuracy.

An important component of BLP's contribution consists of a computationally feasible approach to constructing the moment conditions. As in Berry (1994), the main idea is to invert the non-linear system of market share equations. BLP and Berry suggest nesting this inversion step directly into the parameter search. For complex specifications such as random coefficients, this inversion step may not have an analytic inverse and numerical inversion can be prohibitively slow. BLP propose a contraction-mapping routine to solve this system of equations. This step nests an inner loop contraction mapping into the parameter search. Following the publication of Nevo's (2000b) "A Practitioner's Guide" to implementing BLP, numerous studies have emerged using the BLP approach to estimating discrete choice demand systems with random coefficients.

Our first objective consists of exploring the numerical properties of BLP's contraction mapping approach. The GMM objective function can be called hundreds of times during a numerical optimization over structural parameters; each call to the objective function requires a call to the inner loop. Therefore, it may be tempting to use a less stringent stopping criterion for the inner loop in order to speed up estimation. We show theoretically that any numerical error in the contraction mapping is magnified when considering the numerical error to the overall GMM objective function. Running the inner contraction mapping using a loose stopping criteria propagates numerical error into the GMM objective function, which can cause a smooth optimization routine to stop early and produce parameter estimates that are not a true local minimum. Also, numerical error may prevent the optimization routine from being able to diagnose convergence. The main concern is that researchers may

try to increase the speed of the inner loop by using a looser convergence tolerance. This may lead, unfortunately, to incorrect parameter estimates.

Our second objective consists of proposing a new computational method for implementing the BLP estimator that eliminates the inner loop entirely and, thus, eliminates the potential for numerical inaccuracy discussed above. Following Su and Judd (2007), we recast the BLP problem as a Mathematical Program with Equilibrium Constraints (MPEC). The MPEC method minimizes the GMM objective function subject to a system of nonlinear constraints requiring that the predicted shares from the model equal the observed shares in the data. The minimization of an objective function subject to nonlinear constraints is a standard exercise in nonlinear programming. We prefer the MPEC approach for three reasons. First, there is no numerical error from nested calls, which eliminates the potential for the minimization routine to converge to a point that is not a local minimum of the true GMM objective function, subject to the constraints within a feasibility tolerance, usually set at 10^{-6} . Second, by eliminating the nested calls, the procedure may be faster than the contraction mapping method proposed by BLP. Third, the MPEC algorithm allows the user to relegate all the numerical operations to a single outer loop that can consist of a call to a state-of-the-art optimization package.

BLP is an empirical method and its properties are going to be sensitive to the properties of the data being used. Our third objective is to explore the properties of the data that may cause the nested fixed point (NFP) approach to be slow. We use numerical theory to show that the speed of the NFP contraction mapping is bounded above by a function of what is known as a Lipschitz constant. We derive an analytic expression for the Lipschitz constant in terms of the data and parameter values from the demand model. We then explore which aspects of the data and the data-generating process make the Lipschitz constant higher, and accordingly may make the NFP inner loop slower. In sampling experiments, we find that decreasing the outside good share (by raising the utility intercept) decreases the speed of the NFP estimator. In our sampling experiments, we first compare NFP with a tight stopping criterion for the inner loop to the sloppier approach of using a loose inner loop stopping tolerance. We show that a loose inner loop can lead to parameter estimates that are not true local minima and, depending on the outer loop tolerance, the failure of the optimization routine to report convergence. These numerical findings confirm our theoretical results. The magnitude of the discrepancies in the numerically incorrect parameter estimates from the true parameter values and the numerically correct point estimates is large.

We next directly benchmark MPEC and a correctly-implemented NFP algorithm on fake data where we know the true parameters. We find that NFP with a tight inner loop can be slow when the Lipschitz constant is high. By contrast, MPEC almost always converges and the speed of MPEC appears almost invariant to the Lipschitz constant, which is expected because MPEC does

not nest a contraction mapping. One concern with MPEC may be the large number of parameters in the optimization problem. We increase the number of markets and show that the comparison of the performance of MPEC and NFP does not change as the number of parameters in the optimization problem increases.

It is important to understand that MPEC is statistically the same estimator as BLP. Therefore, the theoretical results on consistency and statistical inference in Berry, Linton and Pakes (2004) apply equally to the contraction mapping and MPEC methods. The related work on identification of the BLP model by Berry and Haile (2008) and Fox and Gandhi (2008) is also agnostic to the actual computational method used in estimation. Our purpose is therefore not to criticize rich structural methods for being too complicated. On the contrary, we view structural methods as a valuable tool in empirical work. Our purpose is to discuss the potential numerical problems that can arise with a complex structural demand model and to offer a practical approach that avoids these problems.

The concerns we raise with the numerical problems from estimating inner loops are magnified as new literatures generalize BLP demand estimators to economically richer models of consumer behavior. As an extension, we consider the discrete choice demand system with forward-looking consumers. We look at the cases, such as in durable goods markets, where consumers can alter a decision to purchase based on expectations about future products and prices (Melnikov 2002, Carranza 2006, Hendel and Nevo 2007, Nair 2007). Gowrisankaran and Rysman (2007) propose the most straightforward extension of the static BLP model. Solving the problem using NFP now involves three numerical loops, adding yet another source of numerical error: the outer optimization routine, the inner inversion of the market share equations, and the inner evaluation of the consumers' value functions (the Bellman equation) for each of the many heterogeneous consumer types. The dynamic programming problem is typically solved with a contraction mapping with the same slow rate of convergence as the BLP market share inversion. Furthermore, Gowrisankaran and Rysman point out that the recursion proposed by BLP may no longer be a contraction mapping for some specifications of dynamic discrete choice models. Hence, the market share inversion is not guaranteed to converge to a solution, which, in turn, implies that the outer optimization routines may not produce the GMM objective function value.

We show that MPEC extends naturally to the case with forward-looking consumers. We optimize the statistical objective function subject to the constraints that Bellman's equation is satisfied at all consumer states and that the market share equations hold. Our approach eliminates both inner loops, thereby reducing these two sources of numerical error. We produce benchmark results that show that MPEC can be faster than NFP under realistic data generating processes. Current research (Dubé, Hitsch and Chintagunta 2008, Lee 2008, Schiraldi 2008) is generalizing BLP to have even more nested

computations than Gowrisankaran and Rysman (2007). The more complicated the model of consumer demand, the greater the advantage of MPEC over traditional inner loop approaches.

Another stream of literature, concerned by the statistical efficiency of GMM estimators, has explored likelihood-based approaches that use additional structure on the joint-distribution of demand and supply (Villas-Boas and Winer 1999; Villas-Boas and Zhao 2005). Jiang et al (2008) propose an alternative Bayesian approach using Markov Chain Monte Carlo methods. In general, likelihood-based approaches still require the numerical inversion of the system of market shares,¹ subjecting them to this additional source of numerical error that MPEC avoids. We outline how one could use MPEC to estimate demand parameters by maximizing the joint likelihood of shares and prices.

Our work on the BLP estimator operates in parallel to Petrin and Train’s (2008) control-function approach, which avoids the inner loop by utilizing additional non-primitive assumptions relating equilibrium prices to the demand shocks. Our proposed MPEC approach also avoids the need for numerical inversion while remaining agnostic about the underlying process (involving the supply side) generating prices and demand shocks – the approach is statistically the same as BLP’s original formulation.

Our assessment of BLP’s numerical properties is broadly related to the recent work by Knittel and Metaxoglou (2008). Knittel and Metaxoglou explore the potential multiple local minima property of the BLP objective function. Our goal is to study the numerical accuracy and speed of finding one local minimum, not to study the broader problems of multiple optima. However, our fake data Monte Carlo experiments routinely find that BLP can recover the true structural parameters using data generated by the model. In a short digression, we examine the same dataset used by Knittel and Metaxoglou and find we cannot replicate their results that the BLP objective function for the cereal data has too many local minima to produce replicable parameter estimates. We choose 50 starting values and find the same local minimum each time, which is the global minimum found by Knittel and Metaxoglou.

To simplify the exposition, hereafter we use “BLP” to refer to the GMM statistical estimator or economic model (random coefficients logit with aggregate demand shocks). We use nested fixed point (NFP) to refer to the traditional algorithm for computing the objective function value, as outlined in BLP (1995). We use MPEC to refer to our alternative, constrained optimization algorithm.

The remainder of the paper is organized as follows. We discuss BLP’s model in Section 2 and their statistical estimator in Section 3. Section 4 provides a theoretical analysis of the numerical properties of BLP’s traditional NFP algorithm. Section 6 presents our alternative MPEC algorithm. Section 7 provides Monte Carlo evidence about the relative performances of the NFP and MPEC algorithms, and considers the estimation error from using a loose stopping tolerance for NFP. The last two sections discuss extensions where MPEC’s advantages over NFP are magnified. First we discuss maximum

¹The transformation-of-variables theorem involves the evaluation of a Jacobian that requires computing the demand shocks numerically.

likelihood estimation, where the need to compute the Jacobian makes MPEC especially useful. Second, we discuss the burgeoning literature on dynamic consumer demand.

2 The Demand Model

In this section, we present the standard random coefficients discrete choice demand model. In most empirical applications, the researcher has access to market shares for each of the available products, but does not have consumer-level information.² The usual modeling solution is to build a system of market shares that is consistent with an underlying population of consumers independently making discrete choices among the various products. The population is in most instances assumed to consist of a continuum of consumers with known mass.

Formally, each market $t = 1, \dots, T$ has a mass M_t of consumers who each choose one of the $j = 1, \dots, J$ products available, or opt not to purchase. Each product j is described by its characteristics $(x_{j,t}, \xi_{j,t}, p_{j,t})$. The vector $x_{j,t}$ consists of K product attributes. The scalar $\xi_{j,t}$ is a vertical characteristic that is observed by the consumers and firms, but is unobserved by the researcher. $\xi_{j,t}$ can be seen as a market and product specific demand shock that is common across all consumers in the market. For each market, we define the J -vector $\xi_t = (\xi_{1,t}, \dots, \xi_{J,t})'$. Finally, we denote the price of product j by $p_{j,t}$.

Consumer i in market t obtains the following indirect utility from purchasing product j

$$u_{i,j,t} = \beta_i^0 + x'_{j,t} \beta_i^x - \beta_i^p p_{j,t} + \xi_{j,t} + \varepsilon_{i,j,t}. \quad (1)$$

The utility of the outside good, or “no-purchase” option, is $u_{i,0,t} = \varepsilon_{i,0,t}$. The consumer i 's preferences consist of the parameter vector β_i , the tastes for each of the K characteristics, and the parameter β_i^p , the marginal utility of income, i 's “price sensitivity”. Finally, $\varepsilon_{i,j,t}$ is an additional idiosyncratic product-specific shock. Let $\varepsilon_{i,t}$ be the vector of all $J + 1$ product-specific shocks for consumer i .

Each consumer is assumed to pick the product j that gives her the highest utility. If tastes, $\beta_i = (\beta_i^0, \beta_i^x, \beta_i^p)$ and $\varepsilon_{i,t}$, are independent draws from the distributions $F_\beta(\beta; \theta)$, with unknown parameters θ , and $F_\varepsilon(\varepsilon)$, respectively, the market share of product j is

$$s_j(x_t, p_t, \xi_t; \theta) = \int_{\{\beta_i, \varepsilon_i | u_{i,j} \geq u_{i,j'} \forall j' \neq j\}} dF_\beta(\beta; \theta) dF_\varepsilon(\varepsilon).$$

To simplify aggregate demand estimation, we follow the convention in the literature and assume ε is distributed type I extreme value, enabling one to integrate it out analytically,

²See Berry, Levinsohn and Pakes (2004) as well as Petrin (2002) for methods incorporating consumer-level data.

$$s_j(x_t, p_t, \xi_t; \theta) = \int_{\beta} \frac{\exp(\beta^0 + x'_{j,t}\beta^x - \beta^p p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^0 + x'_{k,t}\beta^x - \beta^p p_{k,t} + \xi_{k,t})} dF_{\beta}(\beta; \theta). \quad (2)$$

This is the random coefficient logit model.

In BLP, the goal is to estimate the parameters θ characterizing the distribution of consumer random coefficients, $F_{\beta}(\beta; \theta)$. McFadden and Train (2000) prove that a flexible choice of the family $F_{\beta}(\beta; \theta)$ (combined with a polynomial in $x_{j,t}$ and $p_{j,t}$) allows the random coefficient logit model to approximate arbitrarily any vector of choice probabilities (market shares) originating from a random utility model with an observable linear index (meaning no $\xi_{j,t}$ term). Bajari, Fox, Kim and Ryan (2008) prove the nonparametric identification (no finite-dimensional parameter θ) of $F_{\beta}(\beta)$ in the random coefficient logit model without aggregate demand shocks, using data on market shares and product characteristics. Berry and Haile (2008) prove the nonparametric identification of the entire BLP demand model, including allowing for aggregate shocks. Fox and Gandhi (2008) have an alternative identification proof for heterogeneity that can be adapted for market level demand shocks in the same way as Berry and Haile. However, in most applications, more structure is imposed on the family of distributions characterizing $F_{\beta}(\beta; \theta)$ through the choice of the family $F_{\beta}(\beta; \theta)$, with each family member indexed by the estimable finite vector of parameters θ . For example, BLP assume that $F_{\beta}(\beta; \theta)$ is the product of K independent normals, with $\theta = (\mu, \sigma)$, the vectors of means and standard deviations for each component of the K normals.

Typically, the integrals in (2) are evaluated by Monte Carlo simulation. The idea is to generate ns draws of (β, α) from the distribution $F_{\beta}(\beta; \theta)$ and to simulate the integrals as

$$\hat{s}_j(x_t, p_t, \xi_t; \theta) = \frac{1}{ns} \sum_{r=1}^{ns} \frac{\exp(\beta^{0,r} + x'_{j,t}\beta^{x,r} - \beta^{p,r} p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t}\beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})}. \quad (3)$$

In principle, many other numerical methods could be used to evaluate the market-share integrals (Judd 1998, Chapters 7–9).

While a discrete choice model with heterogeneous preferences dates back at least to Hausman and Wise (1978), the inclusion of the aggregate demand shock, $\xi_{j,t}$, was introduced by Berry (1994) and BLP. The demand shock $\xi_{j,t}$ is the natural generalization of demand shocks in the textbook linear supply and demand model. We can see in (2) that without the shock, $\xi_{j,t} = 0 \forall j$, market shares are deterministic functions of the x 's and p 's. In consumer level data applications, the econometric uncertainty is typically assumed to arise from randomness in consumer tastes, ε . This randomness washes out in a model that aggregates over a sufficiently large number of consumer choices (here a continuum). A model without market-level demand shocks will not be able to fit data on market

shares across markets, as the model does not give full support to the data. In the next section, we discuss estimation challenges that arise when $\xi_{j,t}$ is included in the model.

3 The BLP GMM Estimator

We now briefly discuss the GMM estimator typically used to estimate the vector of structural parameters, θ . Like the textbook supply and demand model, the demand shocks, $\xi_{j,t}$, force the researcher to deal with the potential simultaneous determination of price and quantity. To the extent that firms observe $\xi_{j,t}$ and condition on it when they set their prices, the resulting correlation between $p_{j,t}$ and $\xi_{j,t}$ will complicate the estimation of (2). This correlation introduces endogeneity bias.

BLP address the endogeneity of price in demand with a vector of D instrumental variables, $z_{j,t}$. They propose a GMM estimator based on the D moment conditions, $E[\xi_{j,t} | z_{j,t}] = 0$. These instruments can be product-specific cost shifters, although frequently other instruments are used because of data availability. Typically the K non-price characteristics in $x_{j,t}$ are also assumed to be independent of $\xi_{j,t}$ and hence to be valid instruments, although this is not a requirement of the statistical theory. The estimator does not impose a parametric distributional assumption on the demand shocks $\xi_{j,t}$, besides the identifying assumption $E[\xi_{j,t} | z_{j,t}] = 0$.

To form the empirical analog of $E[\xi_{j,t} | z_{j,t}]$ or the often implemented moments $E[\xi_{j,t} z_{j,t}]$, the researcher needs to find the implied values of the demand shocks, $\xi_{j,t}$, corresponding to a guess for θ . The system of market shares, (2), defines a mapping between the vector of demand shocks and the market shares : $S_t = s(x_t, p_t, \xi_t; \theta)$, or $S_t = s(\xi_t; \theta)$ for short. Berry (1994) and Gandhi (2008) prove that s has an inverse, s^{-1} , such that any observed vector of shares can be explained by a unique vector $\xi_t(\theta) = s^{-1}(S_t; \theta)$. For the random coefficients logit specification, we can compute ξ_t using the contraction mapping proposed in BLP. We discuss the properties of the contraction mapping in the next section.

A GMM estimator can now be constructed by using a weighted average of the empirical analog of the moment conditions, $g(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^J \xi_{j,t}(\theta)' z_{j,t}$, where $\xi_t(\theta) = s^{-1}(S_t; \theta)$. For some weight matrix, W , we define the GMM estimator as the vector, θ^{GMM} , that minimizes the function

$$Q(\theta) = g(\theta)' W g(\theta). \tag{4}$$

The statistical efficiency of the GMM estimator can be improved by using other, nonlinear functions of $z_{j,t}$, using an optimal weighting matrix in a second step, or using an efficient one-step method such as continuously-updated GMM or empirical likelihood. However, as we show in the following sections, the numerical precision of the algorithms used to compute $Q(\theta)$ may be equally or more important

from a practical perspective than matters of statistical efficiency.

4 A Theoretical Analysis of the NFP Algorithm

In this section, we theoretically analyze the numerical properties of BLP’s method. The GMM estimator described in section 3 consists of an outer loop to minimize the objective function, $Q(\theta)$, and an inner loop to evaluate this function. Each evaluation of the GMM objective function, $Q(\theta)$, nests a call to a contraction mapping. We call the complete GMM estimator that nests the inner loop the nested fixed point, or NFP, method. Each time the minimization routine calls $Q(\theta)$, the contraction mapping is called T times, once for each market t . If the researcher does not calculate the first and second derivatives of $Q(\theta)$ analytically, many local minimization routines approximate the gradient and Hessian using finite difference methods. The use of finite differences will require many additional calls to $Q(\theta)$ and hence the contraction mapping, proportionately to the dimension of θ .

From a practical perspective, the speed of optimization is determined almost entirely by the number of calls to the contraction mapping and the computation time associated with each run of the inner loop. For these reasons, some practical applications have used a fairly loose convergence criteria to improve speed. In the subsections below, we first provide formal results on the speed of convergence of the inner loop.³ We then show formally how numerical error from the inner loop can propagate into the outer loop, potentially leading to incorrect parameter estimates. One goal of this section is to provide guidelines for researchers in their selection of convergence criteria for the numerical algorithms used to estimate θ^{GMM} . We also theoretically analyze the speed of the NFP algorithm, and discuss when it is likely to be slow.

4.1 The Convergence Rate of the NFP Contraction Mapping

In this section, we derive the rate of convergence of the contraction mapping proposed by BLP to invert the demand system. Recall from section 3 that the evaluation of the GMM criterion, $Q(\theta)$, requires us to evaluate the inverse: $\xi_t(\theta) = s^{-1}(S_t; \theta)$. For a given θ , the inner loop of the NFP estimator solves the share equation for the demand shocks ξ by iterating the contraction mapping

$$\xi_t^{h+1} = \xi_t^h + \log S_t - \log s(\xi_t^h; \theta), \quad t = 1, \dots, T, \quad (5)$$

³Davis (2006) presents an alternative inner-loop method based on a nested optimization problem. It may converge faster than BLP’s contraction mapping.

until the successive iterates ξ_t^{h+1} and ξ_t^h are sufficiently close.⁴ Formally, we choose a small number, for example 10^{-8} or 10^{-10} , for ϵ_{in} as the inner loop tolerance level and require ξ_t^{h+1} and ξ_t^h to satisfy the stopping rule

$$\|\xi_t^h - \xi_t^{h+1}\| \leq \epsilon_{\text{in}} \quad (6)$$

for the iteration $h+1$ where we terminate the contracting mapping (5).⁵ Let $\xi_t(\theta, \epsilon_{\text{in}})$ denote the first ξ_t^{h+1} such that the stopping rule (6) is satisfied. The researcher then uses $\xi_t(\theta, \epsilon_{\text{in}})$ to approximate $\xi_t(\theta)$.

Researchers often find it tempting to loosen the inner loop tolerance if the NFP contraction mapping is slow. Below, we derive formally the theoretical rate of convergence of the inner loop call to the contraction mapping in terms of the economic parameters of the BLP demand model. Numerical theory proves that the convergence of a contraction mapping is linear at best. Linearly convergent algorithms are typically considered to be slow compared to alternative methods, such as Newton's method, for solving nonlinear equations. The numerical performance of a contraction mapping is also sensitive to the stopping tolerance criteria ϵ_{in} . We now state the contraction mapping theorem and discuss how to calculate the linear convergence rate for the inner loop contraction mapping (5) of the BLP estimator.

Theorem 1. *Let $\mathcal{T} : R^n \rightarrow R^n$ be an iteration function and let $S_r = \{\xi \mid \|\xi - \xi^0\| < r\}$ be a ball of radius r around a given starting point $\xi^0 \in R^n$. Assume that \mathcal{T} is a contraction mapping in S_r , meaning*

$$\xi, \tilde{\xi} \in S_r \Rightarrow \|\mathcal{T}(\xi) - \mathcal{T}(\tilde{\xi})\| \leq L \|\xi - \tilde{\xi}\|,$$

where $L < 1$ is called a Lipschitz constant. Then if

$$\|\xi^0 - \mathcal{T}(\xi^0)\| \leq (1 - L)r,$$

the multidimensional equation $\xi = \mathcal{T}(\xi)$ has a unique solution ξ^* in the closure of S_r , $\bar{S}_r = \{\xi \mid \|\xi - \xi_0\| \leq r\}$. This solution can be obtained by the convergent iteration process $\xi^{h+1} = \mathcal{T}(\xi^h)$, for $h = 0, 1, \dots$. The error at the h^{th} iteration is bounded:

$$\|\xi^h - \xi^*\| \leq \|\xi^h - \xi^{h-1}\| \frac{L}{1-L} \leq \|\xi^1 - \xi^0\| \frac{L^h}{1-L}.$$

The Lipschitz constant, L , is a measure of the rate of convergence. At every iteration, the upper

⁴In our implementation of NFP, we iterate over $\exp(\xi)$ to speed up the computation because taking logarithms in MATLAB is slow. However, depending on the magnitude of ξ , the use of the exponentiated form $\exp(\xi)$ in a contraction mapping can lose 3 to 5 digits of accuracy in ξ , and as a result, introduce an additional source of numerical error. For example, if $|\xi_t^h| = -8$ and $|\exp(\xi_t^h) - \exp(\xi_t^{h+1})| = 10^{-10}$, then $|\xi_t^h - \xi_t^{h+1}| = 2.98 \times 10^{-7}$.

⁵ $\|(a_1, \dots, a_b)\|$ is a distance measure, such as $\max(a_1, \dots, a_b)$.

bound for the norm of the error is multiplied by a factor equal to L . A proof of this theorem can be found in many textbooks, such as Dahlquist and Björck (2008). The following theorem shows how a Lipschitz constant for a mapping $\mathcal{T}(x)$ can be expressed in terms of $\nabla\mathcal{T}(x)$, the Jacobian of \mathcal{T} . We then use the Lipschitz constant result to assess an upper bound for the performance of the BLP NFP estimator.

Theorem 2. *Let the function $\mathcal{T}(\xi) : R^n \rightarrow R^n$ be differentiable in a convex set $D \subset R^n$. Then $L = \max_{\xi \in D} \|\nabla\mathcal{T}(\xi)\|$ is a Lipschitz constant for \mathcal{T} .*

The contraction mapping in the BLP estimator is

$$\mathcal{T}(\xi) = \xi + \log S - \log s(\xi; \theta).$$

We define a Lipschitz constant for the BLP contraction mapping \mathcal{T} given structural parameters θ as

$$L(\theta) = \max_{\xi \in D} \|\nabla\mathcal{T}(\xi)\| = \max_{\xi \in D} \|I - \nabla(\log s_j(x_t, p_t, \xi_t; \theta))\|,$$

where

$$\frac{\partial \log(s_j(x_t, p_t, \xi_t; \theta))}{\partial \xi_{jt}} = \begin{cases} \sum_{r=1}^{n_s} \left[\left(\frac{\exp(\beta^{0,r} + x'_{j,t} \beta^{x,r} - \beta^{p,r} p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t} \beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})} \right) - \left(\frac{\exp(\beta^{0,r} + x'_{j,t} \beta^{x,r} - \beta^{p,r} p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t} \beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})} \right)^2 \right], & \text{if } j = l \\ \frac{\sum_{r=1}^{n_s} \frac{\exp(\beta^{0,r} + x'_{j,t} \beta^{x,r} - \beta^{p,r} p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t} \beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})}}{\sum_{r=1}^{n_s} \frac{\exp(\beta^{0,r} + x'_{j,t} \beta^{x,r} - \beta^{p,r} p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t} \beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})}}, & \\ - \sum_{r=1}^{n_s} \left[\left(\frac{\exp(\beta^{0,r} + x'_{j,t} \beta^{x,r} - \beta^{p,r} p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t} \beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})} \right) \left(\frac{\exp(\beta^{0,r} + x'_{l,t} \beta^{x,r} - \beta^{p,r} p_{l,t} + \xi_{l,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t} \beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})} \right) \right], & \text{if } j \neq l \\ \frac{\sum_{r=1}^{n_s} \frac{\exp(\beta^{0,r} + x'_{j,t} \beta^{x,r} - \beta^{p,r} p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t} \beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})}}{\sum_{r=1}^{n_s} \frac{\exp(\beta^{0,r} + x'_{j,t} \beta^{x,r} - \beta^{p,r} p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^{0,r} + x'_{k,t} \beta^{x,r} - \beta^{p,r} p_{k,t} + \xi_{k,t})}}, & \end{cases}$$

For a given vector of structural parameters θ , $L(\theta)$ is the Lipschitz constant for the NFP inner loop. It is difficult to get precise intuition for this expression as it is the norm of a matrix. But, roughly speaking, the Lipschitz constant is related to the matrix of own and cross demand elasticities for the demand shocks, ξ , as the j th element along the main diagonal is $\frac{\partial s_{j,t}}{\partial \xi_{j,t}} \frac{1}{s_{j,t}}$. These expressions are, in turn, related to the degree of asymmetry in the market shares. In section 7.3 below, we use the Lipschitz constant to distinguish between simulated datasets where we expect the contraction mapping to perform relatively slow or fast.

4.2 Determining the Stopping Criteria for the Outer Loop in NFP

This subsection provides guidance on how to select the outer loop tolerance to ensure the outer loop will converge for a given inner loop tolerance. In particular, we show how numerical error from the

inner loop can propagate into the outer loop. We characterize the corresponding numerical inaccuracy in the criterion function, $Q(\theta)$, and its gradient. This analysis then informs the decision of what tolerance to use for the outer-optimization loop to ensure that the optimization routine is able to report convergence. This subsection focuses on ensuring the outer loop can actually converge given the numerical inaccuracy of the inner loop. In a later section, we show how this numerical inaccuracy in $Q(\theta)$ and its gradient can generate numerical inaccuracy in the parameter estimates of θ . In some instances, this inaccuracy could imply that the reported estimates are not a true local minimum of $Q(\theta)$.

Recall that the outer loop of the BLP estimator consists of minimizing the GMM objective function (4). The convergence of this outer loop depends on the choice of an outer loop tolerance level, denoted by ϵ_{out} . In theory, ϵ_{out} should be set to a small number, such as 10^{-5} or 10^{-6} . In practice, we have found cases in the BLP literature where 10^{-2} was used, possibly to offset the slow performance or non-convergence of the minimization routine. As we illustrate in our Monte Carlo simulations below, a loose stopping criterion for the outer loop can cause the routine to terminate early and produce incorrect point estimates. In some instances, these estimates may not even satisfy the first-order conditions for a local minimizer.

We denote by $\xi(\theta, \epsilon_{\text{in}})$ the approximated demand shock corresponding to a given value for θ and an inner-loop tolerance ϵ_{in} that determines the inner-loop stopping rule, (6). We also denote the true demand shock as $\xi(\theta, 0)$. We let $Q(\xi(\theta, \epsilon_{\text{in}}))$ be the programmed GMM objective function with the inner-loop tolerance ϵ_{in} . This more general notation allows us to examine numerical inaccuracy with the programmed inner loop, which is not present in the statistical theory of GMM.

First, we characterize the bias in evaluating the GMM objective function and its gradient at any structural parameters, θ , when there exist inner loop numerical errors. In a duplication of notation, let $Q(\xi)$ be the GMM objective function for an arbitrary guess of ξ . We also use the big-”O” notation whereby $O(T^2)$ is, roughly speaking, a term that grows at the rate T^2 . This notation is a convention in the mathematics literature and is described in many textbooks such as van der Vaart (2000).

Theorem 3. *Let L be the Lipschitz constant for the inner-loop contraction mapping. For any structural parameters θ and with the inner-loop tolerance ϵ_{in} ,*

1. $|Q(\xi(\theta, \epsilon_{\text{in}})) - Q(\xi(\theta, 0))| = O\left(\frac{L(\theta)}{1-L(\theta)}\epsilon_{\text{in}}\right)$
2. $\|\nabla_{\theta}Q(\xi(\theta))|_{\xi=\xi(\theta, \epsilon_{\text{in}})} - \nabla_{\theta}Q(\xi(\theta))|_{\xi=\xi(\theta, 0)}\| = O\left(\frac{L(\theta)}{1-L(\theta)}\epsilon_{\text{in}}\right),$

assuming both $\left\|\frac{\partial Q(\xi)}{\partial \xi}\right\|_{\xi=\xi(\theta, 0)}$ and $\left\|\frac{\partial \nabla_{\theta}Q(\xi(\theta))}{\partial \xi}\right\|_{\xi=\xi(\theta, 0)}$ are bounded.

The notation $O\left(\frac{L(\theta)}{1-L(\theta)}\epsilon_{\text{in}}\right)$ implies that the numerical error in both the objective function and gradient is linear in ϵ_{in} , the tolerance for the inner loop. This result ties back to the fundamental linear

rate of convergence of a contraction mapping. The proof is in the appendix.

Theorem 3 states that the biases in evaluating the GMM objective function and its gradient at any structural parameters are of the same order as the inner-loop tolerance adjusted by the Lipschitz constant for the inner-loop contraction mapping. Recall that a smooth optimization routine converges when the gradient of the objective function is close to zero, by some metric. In the next theorem, we analyze the numerical properties of the gradient. The theorem indicates circumstances in which the outer loop might report convergence despite a numerically inaccurate inner loop.⁶ We also show that the choice of the outer-loop tolerance, ϵ_{out} , should depend on the inner-loop tolerance ϵ_{in} and the Lipschitz constant L . This is important because the outer loop tolerance determines the number of significant digits for the solution. Using a tight outer loop tolerance also helps eliminate spurious local minima.

Theorem 4. *Let $L(\theta)$ be the Lipschitz constant of the inner-loop contraction mapping for a given θ and let ϵ_{in} be the inner-loop tolerance. Let $\hat{\theta} = \arg \max_{\theta} \{Q(\xi(\theta, \epsilon_{\text{in}}))\}$. In order for the outer-loop GMM minimization to converge, the outer-loop tolerance ϵ_{out} should be chosen to satisfy $\epsilon_{\text{out}} = O\left(\frac{L(\hat{\theta})}{1-L(\hat{\theta})}\epsilon_{\text{in}}\right)$, assuming $\left\|\nabla_{\theta}^2 Q(\xi)\Big|_{\xi=\xi(\hat{\theta},0)}\right\|\|\theta - \hat{\theta}\|$ for θ in a neighborhood of $\hat{\theta}$ is bounded.*

The function $\frac{L(\hat{\theta})}{1-L(\hat{\theta})}$ is increasing on $[0, 1]$, the set of valid Lipschitz constants for a contraction mapping. Therefore, if ϵ_{in} is large (the inner loop is loose), then ϵ_{out} must also be large (the outer loop must be loose) for the optimization routine to converge. If the inner loop is slow because L is close to 1, then for a fixed ϵ_{in} , ϵ_{out} should be even larger to ensure convergence. The proof is in the appendix.

An immediate consequence of these results is that the researcher may be tempted to select tolerances based on the convergence of the algorithms, rather than the precision of the estimates themselves. In situations where the inner-loop is slow, a researcher may loosen the inner loop tolerance, ϵ_{in} , to speed convergence of the contraction-mapping. By Theorem 4, the resulting imprecision in the gradient could prevent the optimization routine from detecting a (possibly incorrect) local minimum and converging. In turn, the researcher may be tempted to loosen the outer loop tolerance to ensure convergence of the minimization routine. Besides concerns about imprecision in the estimates, raising ϵ_{out} could also generate an estimate that is not in fact a local minimum.

⁶The numerical error in the gradient convergence test may encourage some researchers to use non-smooth optimization methods. Our experiments with MATLAB's version of a genetic algorithm and the simplex method on the BLP NFP problem suggest that both non-smooth optimizers can report convergence to a point that is not a local minimum, even with a tight inner loop tolerance ϵ_{in} and tight outer-loop tolerance ϵ_{out} . We can verify whether a point is a true local minimum by starting a high-quality smooth optimization routine at that point. If it is a local minimum, the smooth routine will immediately report convergence. For these reasons, we focus on smooth optimizers in this paper.

4.3 Finite Sample Bias in Parameter Estimates from the Inner-Loop Numerical Error

In this section, we discuss the small-sample biases associated with inner-loop numerical error. Assume, given ϵ_{in} , that we have chosen ϵ_{out} to ensure that the algorithm is able to report convergence. Let $\theta^* = \arg \max_{\theta} \{Q(\xi(\theta, 0))\}$ be the maximizer of the finite-sample objective function without numerical error. As economists, now we are interested in the errors in the final estimates, $\hat{\theta} - \theta^*$, from using a loose inner loop.

Theorem 5. *Assume that $\left\| \nabla_{\xi} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right\|$ is bounded and that*

$$O\left(\left\| \xi(\hat{\theta}, \epsilon_{\text{in}}) - \xi(\hat{\theta}, 0) \right\|^2\right) \ll \left| Q(\xi(\hat{\theta}, \epsilon_{\text{in}})) - Q(\xi(\theta^*, 0)) \right| + \left\| \nabla_{\xi} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right\| \left\| \xi(\hat{\theta}, \epsilon_{\text{in}}) - \xi(\hat{\theta}, 0) \right\|.$$

The difference between the finite-sample maximizers with and without numerical error satisfies

$$O\left(\left\| \hat{\theta} - \theta^* \right\|^2\right) \leq \left| Q(\xi(\hat{\theta}, \epsilon_{\text{in}})) - Q(\xi(\theta^*, 0)) \right| + O\left(\frac{L(\hat{\theta})}{1 - L(\hat{\theta})} \epsilon_{\text{in}}\right).$$

The proof is in the appendix. The square of the bias in estimates, $\left\| \hat{\theta} - \theta^* \right\|^2$ is of the same order as $\frac{L(\hat{\theta})}{1 - L(\hat{\theta})} \epsilon_{\text{in}}$, the inner-loop tolerance adjusted by the Lipschitz constant. The “2” in the exponent then implies that the significant digits of the outer loop are only half of that of the inner loop. For example, an inner-loop tolerance of $\epsilon_{\text{in}} = 10^{-6}$ would give the estimated structural parameters in the outer loop three digit accuracy ($\left\| \hat{\theta} - \theta^* \right\| \approx 10^{-3}$). Furthermore, as we have shown in the previous section that the outer loop minimization procedure might not converge if ϵ_{out} is chosen to be 10^{-6} or smaller. As a consequence, accuracy in the estimates and outer loop convergence both require a very tight tolerance criteria for the inner loop, for example, $\epsilon_{\text{in}} = 10^{-10}$.⁷ From a practical perspective, such a high degree of accuracy in the inner loop will slow the speed of convergence of the contraction mapping. In our Monte Carlo experiments, below, we contrast the statistical accuracies of the NFP methods with a loose and tight tolerance for the inner loop.

Note that our theoretical analysis involves local expansions of the objective function and gradients. Consequently, the errors we derive in the parameter estimates are not large. In practice, the errors seen from some sloppy coding techniques may actually be higher than our Taylor series analysis suggests. We document the real-world finite-sample bias in NFP using fake data experiments below.

⁷An even tighter inner loop tolerance is required when finite-difference numerical derivatives are used instead of analytic derivatives.

4.4 Large Sample Bias from the Inner-Loop Numerical Error

The previous section focused only on numerical errors for a finite data set. We now use statistical theory to examine the large-sample properties of the BLP estimator using the NFP algorithm. Before, θ^* was the true minimizer of the finite-sample GMM objective function without any inner-loop numerical errors. Now instead consider θ^0 , the true parameters in the data generating process. Even a researcher with a perfect computer program will not be able to recover θ^0 because of statistical sampling error. Here we explore how numerical errors in the inner loop affect the consistency of the BLP estimator.

Recall that $\hat{\theta}$ corresponds to the minimizer of $Q\left(\xi\left(\hat{\theta}, \epsilon_{\text{in}}\right)\right)$, the biased GMM objective function with the inner-loop tolerance ϵ_{in} . Let $\bar{Q}(\xi(\theta, 0)) = E[Q(\xi(\theta, 0))]$ be the probability limit of $Q(\xi(\theta, 0))$, as either $T \rightarrow \infty$ or $J \rightarrow \infty$, as in Berry, Linton and Pakes (2004). Let $\bar{\theta}$ be the minimizer of $\bar{Q}(\xi(\theta, \epsilon_{\text{in}}))$, the population objective function with the inner-loop tolerance $\epsilon_{\text{in}} > 0$. Clearly, $\theta^0 = \arg \min \bar{Q}(\xi(\theta, 0))$ if the BLP model is identified.

Let asymptotics be in the number of markets, T , and let each market be an iid observation. By standard consistency arguments (Newey and McFadden 1994), θ^* will converge to θ^0 if $Q(\xi(\theta, 0))$ converges to $\bar{Q}(\xi(\theta, 0))$ uniformly, which is usually the case with a standard GMM estimator. Further, the rate of convergence of the estimator without numerical error from the inner loop is the standard parametric rate, \sqrt{T} . By the triangle inequality,

$$\|\hat{\theta} - \theta^0\| \leq \|\hat{\theta} - \theta^*\| + \|\theta^* - \theta^0\| = O\left(\sqrt{\frac{L(\hat{\theta})}{1 - L(\hat{\theta})} \epsilon_{\text{in}}}\right) + O\left(1/\sqrt{T}\right), \quad (7)$$

where $\|\hat{\theta} - \theta^*\| = O\left(\sqrt{\frac{L(\hat{\theta})}{1 - L(\hat{\theta})} \epsilon_{\text{in}}}\right)$ because we showed $O\left(\|\hat{\theta} - \theta^*\|^2\right) \approx O\left(\frac{L(\hat{\theta})}{1 - L(\hat{\theta})} \epsilon_{\text{in}}\right)$ in the previous subsection. These results suggest that the asymptotic bias due to numerical error in the inner loop persists and does not shrink asymptotically. This is intuitive: inner loop error would introduce numerical errors in the parameter estimates even if the population data were used.

4.5 Loose Inner Loop Tolerances and Numerical Derivatives

Most scholars use gradient-based optimization routines, as perhaps they should given that the GMM objective function is smooth. Gradient-based optimization require derivative information, by definition. One approach is to algebraically derive expressions for the derivatives and to manually code them. Our results above assume that the researcher's optimizer has information on the exact derivatives. However, in many applications, such as the dynamic demand model we study below, calculating and coding derivatives can be very time consuming. In this situations, researches may choose to use

numerical derivatives. The gradient is approximated by

$$\nabla_d Q(\xi(\theta, \epsilon_{\text{in}})) = \left\{ \frac{Q(\xi(\theta + de_k, \epsilon_{\text{in}})) - Q(\xi(\theta - de_k, \epsilon_{\text{in}}))}{2d} \right\}_{k=1}^{|\theta|}, \quad (8)$$

where d is a perturbation to an element of θ and e_k is a vector of 0's, except for a 1 in the k th position of e_k . As $d \rightarrow 0$, $\nabla_d Q(\xi(\theta, \epsilon_{\text{in}}))$ converges to $\nabla Q(\xi(\theta, \epsilon_{\text{in}}))$, the numerically accurate derivative of $\nabla Q(\xi(\theta, \epsilon_{\text{in}}))$. However, the ultimate goal of estimation is to minimize the objective function without numerical error. For this, we need the derivatives without numerical error, $\nabla Q(\xi(\theta, 0))$, although that object is not available on the computer.

Lemma 9.1 in Nocedal and Wright (2006) shows that the numerical error in the gradient is bounded,

$$\|\nabla_d Q(\xi(\theta, \epsilon_{\text{in}})) - \nabla Q(\xi(\theta, 0))\|_{\infty} \leq O(d^2) + \frac{1}{d} O\left(\frac{L(\theta)}{1-L(\theta)} \epsilon_{\text{in}}\right).$$

There are two terms in this bound. The $O(d^2)$ term represents the standard error that arises from numerical differentiation, (8). As $d \rightarrow 0$, the $O(d^2)$ term converges to 0. The second term $\frac{1}{d} O\left(\frac{L(\theta)}{1-L(\theta)} \epsilon_{\text{in}}\right)$ arises from the numerical error in the objective function, for a given $\epsilon_{\text{in}} > 0$. The $O\left(\frac{L(\theta)}{1-L(\theta)} \epsilon_{\text{in}}\right)$ term comes from part 1 of Theorem 3. If $\frac{1}{d} O\left(\frac{L(\theta)}{1-L(\theta)} \epsilon_{\text{in}}\right)$ is relatively large, as it is when the inner loop tolerance is loose, then the bound on the error in the gradient is large. In this case, a gradient-based search routine can go completely in the wrong direction, and end up stopping at a parameter far from a local minimum. Therefore, combining loose inner loop tolerances and numerical derivatives will produce an extremely unreliable solver. Note that setting $d \rightarrow 0$ will send the term $\frac{1}{d} O\left(\frac{L(\theta)}{1-L(\theta)} \epsilon_{\text{in}}\right) \rightarrow \infty$. So setting d to be extremely small will only exacerbate the numerical error arising from using a loose inner-loop tolerance.

5 Parameter Errors from Loose Inner-Loop Tolerances in the NFP Algorithm

This section presents evidence that using loose tolerances in the NFP inner loop can lead to incorrect parameter estimates. We show that parameter errors can arise both from fake data and field data. We use the fake data to show that a combination of numerical derivatives and loose inner loop tolerances can lead to grossly incorrect parameter estimates. We use the field data to show that wrong parameter estimates can arise even with closed form derivatives. Both the fake data and real data results are examples: there is no guarantee that any given dataset will combine with a poor implementation of the NFP algorithm to produce incorrect parameter estimates. However, we only need examples in order to show that NFP with loose inner tolerances can produce incorrect parameter estimates.

5.1 NFP Algorithm Implementations

For all NFP implementations, we examine the one-step GMM estimator with Nevo’s (2000) suggestion of using the weighting matrix $W = (Z'Z)^{-1}$, where Z is the $TJ \times D$ matrix of instruments $z_{j,t,k}$.⁸ We use one fake data set and one real data set to show that NFP with loose inner loop tolerances can lead to incorrect parameter estimates.

We use three implementations of NFP for our real data and fake data tests. We use the same data and set of starting values for all three implementations. We use our numerical theory results from section 4 to guide us in the selection of inner and outer loop tolerances for the NFP algorithm. To assess the importance of those findings, we construct three scenarios which we examine for each Monte Carlo experiment. In the first scenario, we explore the implications of a tight outer loop tolerance, set at $\epsilon_{\text{out}} = 10^{-6}$, and a loose inner loop tolerance, set at $\epsilon_{\text{in}} = 10^{-4}$. The former outer loop tolerance is the default setting for most state-of-the-art optimization algorithms. However, from our numerical theory results, we know the latter inner loop tolerance is too large. One could think of this scenario as representing the “frustrated researcher” who loosens the inner loop to speed the apparent rate of convergence. In the second scenario, we explore the results from Theorem 4, whereby the loose inner loop tolerance could, in turn, prevent the outer loop from converging. Specifically, we keep $\epsilon_{\text{in}} = 10^{-4}$ and set $\epsilon_{\text{out}} = 10^{-2}$. One can think of this scenario as representing the attempt of the researcher to loosen the outer loop to force it to converge, even though in practice the converged point may not actually satisfy the first-order conditions. In our third scenario, we implement the “best practice” settings for the NFP algorithm with $\epsilon_{\text{in}} = 10^{-14}$ and $\epsilon_{\text{out}} = 10^{-6}$.

For all implementations of NFP, we use the same programming environment (MATLAB) and the same optimization package (KNITRO using the TOMLAB interface). We selected MATLAB because this is a commonly-used software package among practitioners. We also selected the KNITRO optimization package instead of MATLAB’s built-in optimization routines as the former is a highly-respected, state-of-the-art solver in the optimization community (Byrd, Nocedal and Waltz 1999). For our fake data example, we use numerical derivatives. For our real data example, we also supply derivatives for each algorithm because all local optimization methods improve if the user supplies exact derivatives of the objective function.⁹

We also customized several aspects of the NFP algorithm to increase speed. In the case of NFP, the

⁸We choose a simple weighting matrix because our focus is on comparing algorithms, not finding the most statistically efficient estimator.

⁹Another option is to use automatic differentiation software. Automatic differentiation software is automatically used by some languages, such as AMPL, and can be accessed with the TOMLAB interface for MATLAB. Our experience has been that automatic differentiation is very slow for NFP. Also, software packages like AMPL are impractical for NFP algorithms because AMPL is a problem definition language, not a general purpose programming language like MATLAB. Therefore, we use MATLAB for all our empirical analysis. However, in practice, many users may find AMPL more convenient for the MPEC implementation. One warning: the automatic differentiation overhead in AMPL uses lots of computer memory.

most notable speed improvements came from exploiting as much as possible the built-in linear algebra operations (“vectorization”) for the inner loop. In addition, we exploited the normality assumption for $F_\beta(\beta; \theta)$ to concentrate the means out of the parameter search under the NFP algorithm, as suggested in Nevo (2000b). Therefore, the NFP algorithm can be recast to search only over the standard deviation of the random coefficients, rather than both the means and standard deviations. Relaxing the normality assumption would prevent the use of this simplification (except perhaps in other location and scale families), which could improve the relative speed performance of MPEC over NFP even further.

The Fake Data Generating Process

We use the demand model from section 2. In this section we describe a data generating process for a base case. The individual experiments perturb aspects of the data generating process from this base case. We allow for $K = 3$ observed characteristics, in addition to prices. We also estimate a random coefficient on the intercept, β_i^0 , which models the relative attractiveness of purchasing any of the products instead of the outside good. β_i^p , the price coefficient, is also random.

We focus on markets with a fairly large number of products, $J = 75$, to ensure that our results are not due to sampling error. We also consider an intermediate number of statistically independent markets, here $T = 25$. Although not reported, we noticed large biases in the mean and standard deviation of the intercept, β_i^0 , as well as functions of the parameters (like price elasticities) when a small number of markets was used. Intuitively, the moments of β_i^0 are identified in part from the share of the outside good, and more markets are needed to observe more variation in the outside good’s shares.

For product j in market t , let

$$\begin{bmatrix} x_{1,j,t} \\ x_{2,j,t} \\ x_{3,j,t} \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & -0.8 & 0.3 \\ -0.8 & 1 & 0.3 \\ 0.3 & 0.3 & 1 \end{bmatrix} \right).$$

Likewise, $\xi_{j,t} \sim N(0, \sigma_\xi^2)$, with the default $\sigma_\xi^2 = 1$. Price is

$$p_{j,t} = |0.5 \cdot \xi_{j,t} + e_{j,t}| + 1.1 \cdot \left| \sum_{k=1}^3 x_{k,j,t} \right|,$$

where $e_{j,t} \sim N(0, 1)$ is an innovation that enters only price. Prices are always positive. Prices are endogenous as $\xi_{j,t}$ enters price. For each product j in market t , there is a separate vector $z_{j,t}$ of $D = 6$ instruments. A powerful instrument must be correlated with $p_{j,t}$ and a valid instrument must not be correlated with $\xi_{j,t}$. Each instrument $z_{j,t,d} \sim N(\frac{1}{4}p_{j,t}, 1)$.

The goal is to estimate the parameter vector θ in $F_\beta(\beta; \theta)$, the distribution of the random coeffi-

cients. To maintain consistency with the application in BLP (1995) and the related empirical literature, we assume independent normal random coefficients on each product characteristic and the intercept. Thus, $F_{\beta}(\beta; \theta)$ is the product of five independent normal distributions ($K = 3$ attributes, price and the intercept) characterized by means and standard deviations contained in θ . The true values of the moments of the random coefficients $\beta_i = \{\beta_i^0, \beta_i^1, \beta_i^2, \beta_i^3, \beta_i^p\}$ are $E[\beta_i] = \{-1, 1.5, 1.5, 0.5, -3.0\}$ and $\text{Var}[\beta_i] = \{0.5, 0.5, 0.5, 0.5, 0.2\}$.

Our focus is not on numerical integration error, so we use the same set of 100 draws to compute market shares in the data generating and estimation phases using (3). By using the same draws, we have no simulation error. By using 100 draws, our code is proportionately faster, allowing us to try more starting values and, in a later section, more Monte Carlo replications. The computational times reported in a later section should be inflated by the ratio of the number of draws that would be used in real data applications (say 4000) to the number we use, 100.

5.2 The Nevo (2000b) Cereal Data and Knittel and Metaxoglou (2008)

We use the cereal data set from Nevo (2000b) to assess whether NFP with loose inner loop tolerances can produce incorrect parameter estimates with real data. We refer the reader to Nevo (2000b) for a description of these data.

We first verify the claim of Knittel and Metaxoglou (2008) regarding the reliability of the BLP GMM estimator applied to these same data. They report that the parameter estimates are extremely sensitive to the starting values used for the NFP algorithm because many true local minima exist in the GMM objective function. We agree that the BLP problem is not convex and, therefore, may potentially generate multiple optima. However, we find the extremity of the problems reported in Knittel and Metaxoglou surprising in light of our own experiments. Careful inspection of the output provided by Knittel and Metaxoglou on their website reveals that those runs that did not produce the lowest objective function value were typically cases for which the outer loop optimization had failed to converge: these were not local optima.¹⁰ Therefore, we use our NFP code with 50 different starting points to estimate the demand model with Nevo’s cereal data set. We set the inner loop tolerance to be 10^{-14} . The starting values consist of 50 draws from the standard normal distribution, which is the same choice made by Knittel and Metaxoglou.¹¹ For each of the 50 runs, our NFP code finds the same objective function value, 4.5615, which is also the lowest objective value found by Knittel and Metaxoglou (2008).¹² We therefore respectfully disagree with Knittel and Metaxoglou’s

¹⁰The Knittel and Metaxoglou website is http://www.econ.ucdavis.edu/faculty/knittel/KM_website.html. We refer here to results reported for the October 1, 2008 version of Knittel and Metaxoglou (2008).

¹¹We also experimented with multiplying the starting values by the solution reported in Nevo (2000b). The results were similar.

¹²We also attempt to replicate the experiment in Knittel and Metaxoglou (2008) by using the MATLAB code by Nevo (2000b), except that we use the KNITRO solver as the search algorithm. For 50 starting points, KNITRO only

interpretation of their findings. With multiple starting values, careful implementation of the numerical procedures, and state-of-the-art optimization solvers, the BLP GMM estimator appears to produce reliable estimates.¹³

5.3 Fake Data, Numerical Derivatives and False Parameter Estimates

For NFP, the numerical theory in section 4 raises several concerns about the common practice of setting the tolerance, ϵ_{in} , too high (too loose). Section 4.5 shows that a combination of a loose inner loop, numerical derivatives and a smooth optimization routine can produce incorrect parameter estimates. Also recall that Theorem 4 shows that if ϵ_{in} is too loose, ϵ_{out} must be set to be too loose in order for the routine to be able to report convergence.

In this subsection, we explore empirically the problems with loose inner loop tolerances and numerical derivatives. We create one simulated fake dataset, using the data generating process from section 5.1. Holding the simulated data fixed, we first compare the estimates produced from 100 randomly-chosen starting values for the own-price demand elasticities. We run each of the three NFP implementations described in section 5.1 for each of the 100 vectors of starting values. Table 1 reports the results for the 100 different starting values. The first row reports the fraction of runs for which the routine reports convergence. As Theorem 4 shows, if the inner loop tolerance is a loose $\epsilon_{\text{in}} = 10^{-4}$ and the outer loop tolerance a standard value of $\epsilon_{\text{out}} = 10^{-6}$, the routine will never report convergence. Column one confirms this finding as only 2% of the runs with a loose inner loop and tight outer loop converge. In contrast, column two indicates that the algorithm is more likely to converge (30% of the runs) when we also loosen the tolerance on the outer loop. As we will show below, this semblance of convergence is merely an artifact of numerical imprecision that leads to misleading estimates. Finally, NFP with tight tolerances converges in 95% of the runs.

To diagnose the quality of the estimates, the second row of Table 1 shows the fraction of runs where the reported GMM objective function value was within 1% of the lowest objective function that we numerically found across all three NFP implementations and all 100 starting values (300 cases). We call this value the “global” minimum, although of course we cannot prove we have found the true global minimum. In the first two columns, corresponding to the scenario with a loose inner loop and the scenario with a loose inner and outer loop respectively, none of the 100 starting values lead to finding the global minimum.¹⁴ In contrast, NFP tight found the global minimum for 25% of the time.

converges for 25 of the 50 runs. However, all 25 successful runs converge to the same solution with objective value 4.5615.

¹³We also use MATLAB’s genetic algorithm routine for one run. The genetic algorithm finds a point with the objective function value 98.4270, which is clearly an order of magnitude higher than 4.5615, the minimum we find using the gradient-based method. We then start KNITRO from this point found by the generic algorithm and KNITRO finds the solution with objective value 4.5615. So the genetic algorithm does not always find a local minimum.

¹⁴Even with loose inner loop tolerances, the GMM objective value function is accurate to a few decimal places. The

Table 1: Three NFP Implementations: Varying Starting Values for One Fake Dataset, with Numerical Derivatives

	NFP Loose Inner	NFP Loose Both	NFP Tight	Truth
Fraction Reported Convergence	0.02	0.30	0.95	
Frac. Obj. Fun. < 1% Greater than “Global” Min.	0.0	0.0	0.25	
Mean Own Price Elasticity Across All Runs	-12.28	-12.30	-5.77	-5.68
Std. Dev. Own Price Elasticity Across All Runs	19.44	19.43	0.0441	
Lowest Objective Function Value	0.0217	0.0327	0.0169	
Elasticity for Run with Lowest Obj. Value	-5.89	-5.63	-5.77	-5.68

We used 100 starting values. The NFP loose inner loop implementation has $\epsilon_{\text{in}} = 10^{-4}$ and $\epsilon_{\text{out}} = 10^{-6}$. The NFP loose both implementation has $\epsilon_{\text{in}} = 10^{-4}$ and $\epsilon_{\text{out}} = 10^{-2}$. The NFP tight implementation has $\epsilon_{\text{in}} = 10^{-14}$ and $\epsilon_{\text{out}} = 10^{-6}$. We use numerical derivatives using KNITRO’s built-in procedures.

NFP tight should not find the global minimum every time, because a gradient-based optimization routine may indeed converge to a local minimum.

The third and fourth rows of Table 1 provide measures to assess the economic implications of our different implementations. We use estimated price elasticities to show how naive implementations could produce misleading economic predictions. In the third row, we report the mean own price elasticity, across all 100 starting values, all $J = 25$ products and all $T = 75$ markets:

$$\frac{1}{100} \sum_{h=1}^H \frac{1}{T} \sum_{t=1}^T \frac{1}{J} \sum_{j=1}^J \eta_{j,t}^p(\hat{\theta}^h),$$

where $\hat{\theta}^h$ is the vector of parameter estimates for the h th starting value and $\eta_{j,t}^p(\hat{\theta}^h)$ is the own price elasticity of firm j in market t , at those parameters. The fourth row reports the standard deviation of the mean own price elasticity across all 100 runs: $\frac{1}{T} \sum_{t=1}^T \frac{1}{J} \sum_{j=1}^J \eta_{j,t}^p(\hat{\theta}^h)$.

Beginning with the third row, first note that in the final column we report the own-price demand elasticity evaluated at the true parameter values: -5.68. As we hoped, NFP with a tight tolerance produces an estimate near the truth, -5.77. Also, even though only 25% of estimates converged to the “global minimum”, the other local minima produce very similar own-price demand elasticities: the standard deviation across starting values is only 0.0441. On the other hand, we immediately see that the loose tolerance implementations of NFP produce mean elasticities that are not nearly as close to the truth as NFP with a tighter inner loop tolerance. The mean of the NFP loose inner implementation is -12.28, nearly twice in absolute value the true value of -5.68. The loose both results are nearly identically. The standard deviations of own price elasticities for the loose inner loop tolerances are huge: 19.4. With a loose inner loop tolerance and numerical derivatives, section 4 shows that there is no reason to expect the NFP algorithm to produce correct parameter estimates.

values reported in the fifth row of Table 1 are identical whether the reported parameter estimates are evaluated at a NFP objective function with $\epsilon_{\text{in}} = 10^{-4}$ or $\epsilon_{\text{in}} = 10^{-14}$.

Table 2: Three NFP Implementations: Varying Starting Values for Nevo’s Cereal Dataset, with Closed-Form Derivatives

	NFP Loose Inner	NFP Loose Both	NFP Tight
Fraction Reported Convergence	0.0	0.81	1.00
Frac. Obj. Fun. < 1% Greater than “Global” Min.	0.0	0.0	1.00
Mean Own Price Elasticity Across All Runs	-3.75	-3.69	-7.43
Std. Dev. Own Price Elasticity Across All Runs	0.03	0.08	~0
Lowest Objective Function Value	15.3816	15.4107	4.5615
Elasticity for Run with Lowest Obj. Value	-3.77	-3.77	-7.43

We use the same 25 starting values for each implementation. The NFP loose inner loop implementation has $\epsilon_{\text{in}} = 10^{-4}$ and $\epsilon_{\text{out}} = 10^{-6}$. The NFP loose both implementation has $\epsilon_{\text{in}} = 10^{-4}$ and $\epsilon_{\text{out}} = 10^{-2}$. The NFP tight implementation has $\epsilon_{\text{in}} = 10^{-14}$ and $\epsilon_{\text{out}} = 10^{-6}$. We manually code closed-form derivatives.

One question is whether a researcher who tried 100 starting values could get close to the true estimates. If “close” is defined by getting one significant digit of the mean own-price elasticity correct, the answer for this particular dataset is “yes”. The fifth row reports the lowest found objective function values found across all 100 starting values by the three algorithms: NFP with a tight tolerance finds an objective function that is lower than the two loose implementations. The reported elasticities for NFP with loose inner loops are -5.89 and -5.62, compared to the numerically correct -5.77 from NFP tight. What is happening is that the NFP implementations with the loose inner loop tolerances tend to stop near the starting values. By using 100 starting values, the researcher is exploring 100 regions of the objective function. It is equivalent to just evaluating the objective function at 100 points and taking the final estimates to be the minimum. However, there is no guarantee that 100 starting values will lead to “close” estimates in other datasets. Indeed, in the next we show an example where even the elasticities corresponding to the lowest objective function value have even more numerical error in them.

5.4 Parameter Errors with the Cereal Data and Closed-Form Derivatives

There are at least three concerns one might have with the previous section’s fake data example. First, perhaps real data does not have the problems we found. Second, the example relied on numerical derivatives; perhaps coding closed-form derivatives eliminates all concerns with achieving incorrect parameter estimates with NFP with loose inner loop errors. Third, the incorrect elasticity estimates in Table 1 were really variable across starting values. A researcher who tried even a few starting values and found the wildly different elasticity estimates would diagnose that something is wrong. A careful researcher might then explore settings like the inner loop tolerance, and eventually fix the implementation error. This section uses Nevo’s cereal data to produce an example of incorrect parameter estimates that is robust to these concerns.

The results in Table 2 are of the same format as Table 1. As Theorem 4 predicts, in row 1 we find that 0% of the NFP loose inner loop starting values converge. Loosening the outer loop is one approach to finding convergence; the second column finds that 81% of starting values report convergence when this is done. 100% of the starting values converge for NFP tight. The second row shows that 100% of the NFP tight starting values find the global minimum, 4.5615, in Nevo’s cereal data. None of the NFP loose tolerance implementations find the global minimum.

The loose inner loop and loose both methods find a mean own-price elasticity of -3.75 and -3.69, respectively. This is about half the value of -7.43 found with NFP tight. Further, the estimates are all tightly clustered around the same points. With standard deviations of 0.03 and 0.08 for the loose inner loop methods, the answers are consistently wrong across runs. The fifth row shows the smallest objective function values found by the loose inner loop and loose both routines are 15.38 and 15.41, respectively. These are far from the true “global” minimum of 4.56.

These results show that a naive but otherwise careful researcher might feel that his or her estimates were correct because even trying 25 different starting values always produce around the same estimates. Even if the researcher correctly coded the derivatives in closed form and used a high-quality, professional optimizer like KNITRO, the NFP loose inner and loose both implementations can consistently converge to the wrong elasticity, and the elasticity can be half of the true value. Thus, there is no diagnostic that a researcher can do that will detect all types of numerical error. With Nevo’s cereal dataset, an inner loop tolerance that is too loose will lead to consistent but consistently wrong own-price elasticity estimates. Only using an a priori theoretically correct setting, like a tight inner loop tolerance, will avoid these errors.

6 A Constrained Optimization Approach to Improve Speed

We have established that only NFP with a tight inner loop tolerance can produce reliable parameter estimates. According to Theorem 5, if we wish to achieve the default numerical precision in the outer loop of 10^{-6} , we need to set the NFP inner loop tolerance to 10^{-12} or tighter, for reliable parameter estimates. Using a tight inner loop means NFP may be slow. Further, in the previous section, we established that the NFP method’s inner loop converges linearly and can be slow when the Lipschitz constant is close to 1. A slow inner loop might cause researchers to choose loose tolerances for the inner loop, which might lead to problems in establishing the convergence of the outer loop as well as errors in the reported parameter estimates.¹⁵

¹⁵Alternative methods to a contraction mapping for solving systems of nonlinear equations with faster rates of convergence typically have other limitations. For instance, the traditional Newton’s method is only guaranteed to converge if the starting values are close to a solution, unless one includes line-search or trust-region procedure subject to some technical assumptions. In general, most practitioners would be daunted by the task of nesting a hybrid Newton

In this section, we propose an alternative algorithm based on Su and Judd’s (2007) constrained optimization approach for estimating structural models. Below we show that the MPEC approach generates the same solution as NFP. MPEC can save computation time while completely avoiding issues of numerical precision by eliminating the inner loop of the NFP algorithm. In their original paper, Su and Judd focus more on solving for the unknown variables in economic models, such as value functions in single-agent dynamic programming problems and the entry probabilities of rival firms in static Bertrand entry games with multiple equilibria. We apply this insight to the recovery of the unobserved demand shocks that enter the criterion function during estimation of a structural model. In particular, we present a constrained optimization formulation for random-coefficients demand estimation.

If W is the GMM weighting matrix, our constrained optimization formulation is

$$\begin{aligned} \min_{\theta, \xi} \quad & g(\xi)' W g(\xi) \\ \text{subject to} \quad & s(\xi; \theta) = S \end{aligned} \tag{9}$$

The moment condition term $g(\xi)$ is just

$$g(\xi) = \frac{1}{T} \sum_{t=1}^T \xi_t z_t.$$

In MPEC, the market share equations are introduced as nonlinear constraints to the optimization problem. The objective function is specified primitively as a function of the demand shocks ξ . The main difference compared to the traditional NFP method is that we optimize over both the aggregate demand shocks ξ and the structural parameters θ . We do not use NFP’s inner loop to enforce $\xi = \xi(\theta)$ for every guess of θ ; rather we impose that the predicted shares equal the actual shares in the data only at the solution to the minimization problem.

The next theorem shows the equivalence of the first-order conditions between the NFP method (4) and the constrained optimization formulation (9). Hence, any first-order stationary point of (4) is also a stationary point of (9), and vice versa.

Theorem 6. *The set of first order conditions to the MPEC minimization problem in (9) is equivalent to the set of first order conditions to the true (no numerical error) GMM inner loop / outer loop method that minimizes (4).*

The main benefit of the MPEC formulation is that it circumvents the need for the inner loop. By eliminating the inner loop, MPEC is less prone to numerical errors and is potentially faster. We discuss these benefits below.

method customized to a specific demand problem inside the outer optimization over structural parameters.

The constrained optimization defined by (9) can be solved using modern nonlinear optimization solvers developed by researchers in numerical optimization. Unlike the NFP algorithm, where users need to exercise caution in the choice of tolerance levels for both inner and outer loops, the defaults on feasibility and optimality tolerances in nonlinear optimization solvers for constrained optimization are usually sufficient. These default tolerances have been established to work well in hundreds or thousands of papers in the numerical analysis literature. The default tolerances are usually sufficient because the market share equations and GMM objective function (without an inner loop) are exposed to the optimization routine. In short, MPEC lets a state-of-the-art optimization algorithm handle all of the computational aspects of the problem. In contrast, with NFP, the researcher needs to customize a nested-fixed-point calculation, which could result in naive errors.

In addition to simplifying implementation, bypassing the inner loop reduces several sources of numerical error that could, possibly, lead to non-convergence. We have detected some common practices with the coding of the inner loop that could naively lead to numerical error. These include loose choices for the inner loop tolerance (as discussed previously) and an adjustable inner-loop tolerance that is loosened for parameter values deemed “far” from the solution to the outer loop.¹⁶ MPEC relegates all numerical calculations to a single call to the outer-loop, which is solved using a state-of-the-art optimization package, rather than the user’s own customized algorithm.

Our approach can also create substantial speed advantages. As we showed in the previous section, the contraction mapping in the NFP algorithm might be slow as the Lipschitz constant approaches one. By contrast, the MPEC method does not nest any contraction mappings and so we expect its speed to be relatively invariant to the Lipschitz constant. Most optimization solvers for smooth problems use Newton-type methods to solve the Karush-Kuhn-Tucker system of the first-order optimality conditions. Newton’s method is quadratically convergent, faster than the linear rate of the contraction mapping that is the NFP inner loop. Another potential source of acceleration of speed comes from the fact that our approach allows constraints to be violated during the solving process. In contrast, the NFP algorithm requires solving the share equation (2) *exactly* for every parameter θ examined in the outer, optimization loop. Modern numerical optimization solvers do not enforce that the constraints are satisfied at every iteration; it suffices that the constraints hold at the solution. This flexibility avoids wasting computational time on iterates away from the true parameters. Still another potential speed advantage is that the outer algorithm has more information: the optimization

¹⁶The trick consists of using a loose inner loop tolerance when the parameter estimates appear “far” from the solution and switching to a tighter inner loop tolerance when the parameter estimates are “close” to the solution. The switch between the loose and tight inner loop tolerances is usually based on the difference between the successive parameter iterates, e.g, if $\|\theta^{k+1} - \theta^k\| \leq 0.01$, then $\epsilon_{\text{in}} = 10^{-8}$; otherwise, $\epsilon_{\text{in}} = 10^{-6}$. Suppose that the following sequence of iterates occur: $\|\theta^{k+1} - \theta^k\| \geq 0.01$ ($\epsilon_{\text{in}} = 10^{-6}$), $\|\theta^{k+2} - \theta^{k+1}\| \leq 0.01$ ($\epsilon_{\text{in}} = 10^{-8}$), and $\|\theta^{k+2} - \theta^{k+1}\| \geq 0.01$ ($\epsilon_{\text{in}} = 10^{-6}$). The NFP objective value can oscillate because of the use of two difference inner loop tolerances. This oscillation can prevent the NFP approach from converging.

routine is exposed to the constraints, the derivatives of the constraints and of the objective function, and the sparsity pattern of the constraints. On sparsity, recall that demand shocks for market t do not enter the constraints for market $t + 1$. Therefore, this constrained optimization problem is highly sparse.

Most constrained optimization solvers are based on sequential quadratic programming or interior point methods. As stated earlier, these solvers use Newton-based methods. Economists are often skeptical about Newton’s method because it might not converge if the starting point is far away from the solution. While this perception is true for the purest textbook version of Newton’s method, modern Newton-like methods incorporate a line-search or a trust-region strategy to give more robustness to the choice of starting values. We refer readers to Nocedal and Wright (2006) and Kelley (1995, 1999, 2003) for further details on modern optimization methods for smooth objectives and constraints.

Finally, our implementation of MPEC for the BLP model is slightly more sophisticated than the simple explanation in (9). We actually treat the moments as separate parameters, so that the problem being solved is

$$\begin{aligned} \min_{\theta, \xi, \eta} \quad & \eta' W \eta \\ \text{subject to} \quad & g(\xi) = \eta \quad . \\ & s(\xi; \theta) = S \end{aligned} \tag{10}$$

The solution to this new problem is the same as (9). The objective function is now a simple quadratic, $\eta' W \eta$, rather than a more complex, direct function of ξ ; the additional constraint $g(\xi) - \eta = 0$ is linear in both ξ and η and, hence, does not add additional difficulties to the original problem. Computationally, the advantage with this equivalent formation is that we increase the sparsity of the constraint Jacobian and the Hessian of the Lagrangian function by adding the additional variables and constraints. In numerical optimization, it is often easier to solve a large but sparse problem than a small but dense problem. Another advantage of MPEC over NFP is that the objective function and constraints in MPEC are likely more “smooth” or less “nonlinear” in the unknowns than the NFP objective function is in θ . In NFP, the mapping from θ to the objective function value uses the very nonlinear inner loop transformation, while no such inner loops are used by MPEC. Thus, MPEC may be a “smoother” nonlinear programming problem.

A common response to practitioners when they first hear about MPEC is that maximizing over a large number of parameters is a numerically daunting challenge. Below we show that this need not be the case. Indeed, the performance comparison of MPEC and NFP may be relatively constant as the number of products and markets increases.

7 Speed Comparisons of MPEC and NFP

NFP with a tight inner loop will produce correct parameter estimates if many starting values are used. However, NFP can be slow on some datasets. This section uses fake data and the Nevo cereal to compare the speed of MPEC and NFP. We present examples where MPEC performs better than NFP. This is not meant to be a theorem: there could be cases where NFP is faster than MPEC. We now show that, in many situations, NFP may be computationally impractical in terms of execution time. In contrast, we will show that MPEC’s execution time appears to be relatively invariant across these situations. Our approach exploits the Lipschitz constant for the BLP contraction mapping derived in section 4.1. We conjecture that data with a higher Lipschitz constant, and hence a higher upper bound on the rate of convergence of the inner loop, may slow NFP estimation. The idea will be to manipulate various components of the data-generating process in order to measure their respective impact on the Lipschitz constant. We have no reason to believe cases exist where MPEC grows really slow with some equivalent of a Lipschitz constant. Therefore, we suspect that MPEC will be more robust against extremely slow performance. Keep in mind that in these slow-performing cases where a researcher will be tempted to loosen the inner loop tolerance, leading to the problem of incorrect parameter estimates that we earlier highlighted.

7.1 NFP and MPEC Implementations

We code NFP and MPEC using closed-form derivatives. As the proof of Theorem (6) shows, the components of these derivatives are the same for both methods. We use the quadratic form of MPEC in (10). We give the sparsity pattern of the constraints to the optimization routine, for MPEC.

An important point for our speed comparison is the choice of starting values. We always use five starting values. For all implementations of the NFP algorithm, for our first starting value we use $\frac{1}{2} |\beta^{2SLS}|$ for each parameter as starting values for the normal standard deviations. Here β^{2SLS} is the vector of coefficient estimates from the logit model without random coefficients, which can be estimated using linear instrumental variable methods. Our starting values are based on a simple conjecture that standard deviations tend to be lower than means (in absolute value), as has been roughly the case in our empirical experience. Our choice of starting values ensures that the parameters are about the correct magnitudes.¹⁷ For other starting values, we multiply $\frac{1}{2} |\beta^{2SLS}|$ by a vector of random numbers, each element of which is distributed as a uniform with support $[0, 3]$. This choice of support keeps the dispersion parameters positive, but ensures that we look over a relatively wide range of values.

For MPEC, we effectively use the same starting values as we do for NFP; we pick values such that

¹⁷Not using 2SLS to inform the starting values will of course lead to longer run times.

the two algorithms are initialized to have the same objective function value.¹⁸ For each NFP starting value, we run the inner loop once and use this vector of demand shocks and mean taste parameters as starting values for MPEC. This is our attempt to equalize the starting values across NFP and MPEC.¹⁹

As before, we use 100 simulation draws for both MPEC and NFP. For the fake data experiments, the same 100 simulation draws to generate the data and to estimate the model. This shuts down simulation error. Raising the number of simulation draws to a more reasonable number, say 10,000, would increase the CPU times of both MPEC and NFP by about 100 times. So the reported times below are 100 times too slow, compared to an actual empirical investigation.

7.2 Base Fake Data Case

Here we define a base fake data case, which is then perturbed to vary the Lipschitz constants in the examples that follow. The model is nearly the same as Section (5.1). We use $T = 50$ to speed the runs somewhat. The mean of the random coefficients is $E[\beta_i] = \{-1, 1.5, 1.5, 0.5, -3.0\}$. The prices are $p_{j,t} = |\xi + \tilde{p}_{j,t} \cdot 1.5|$, where $\tilde{p}_{j,t} = 0.5 + \eta_{j,t} \cdot 3.5 + 0.5 \cdot \sum_{k=1}^3 x_{k,j,t}$ and $\eta_{j,t}$ is a uniform(0,1) random variable. Likewise, $z_{j,t,d} = \tilde{\eta}_{j,t,d} + \frac{1}{4}\tilde{p}_{j,t}$, where $\tilde{\eta}_{j,t,d}$ is another uniform(0,1) random variable. For each table below, we calculate 20 or 30 different fake data sets, and reported means across these 20 or 30 replications.

7.3 Lipschitz Constants

Recall that the Lipschitz constant derived in section 4.1 is related to the demand sensitivity to the unobserved quality, $\xi_{j,t}$. Moreover, this demand sensitivity is roughly related to the degree of asymmetry in market shares. Therefore, we experiment with different features of the data-generating process that affect the degree of share asymmetry. Table 3 reports the Lipschitz constant for the base-case data-generating process of section 5.1. Each cell reports the mean of the Lipschitz constant evaluated at the true parameter values across 30 data sets / replications.

In our first experiment, reported in the first column of Table 3, we manipulate the scale of the parameters, β_i . We multiply the β_i of each of our ns simulated consumers in the data generating process by one of the constants listed in the table. We find that the Lipschitz constant is non-monotone in the scale, with the constant first falling and then rising again. This non-monotonicity comes from the fact that our manipulation also changes the levels of the market shares. Nevertheless, holding the

¹⁸Our MATLAB code is efficiently parallelized across multiple cores. However, we report CPU times and not clock times.

¹⁹Adding the NFP inner loop takes two lines of code once MPEC has been coded, so it is not unreasonable to expect a practitioner to be able to reproduce our choice of MPEC starting values.

Table 3: Lipschitz Constants for the NFP Algorithm

Parameter Scale		Std. Dev. of Shocks ξ		# of Markets T		Mean of Intercept $E[\beta_i^0]$	
Altered Value	Mean Lipschitz	Altered Value	Mean Lipschitz	Altered Value	Mean Lipschitz	Altered Value	Mean Lipschitz
0.01	0.985	0.1	0.808	25	0.860	-2	0.771
0.1	0.971	0.25	0.813	50	0.871	-1	0.871
0.50	0.887	0.5	0.832	100	0.888	0	0.936
0.75	0.865	1	0.871	200	0.888	1	0.971
1	0.871	2	0.934			2	0.988
1.5	0.911	5	0.972			3	0.996
2	0.938	20	0.984			4	0.998
3	0.970						
5	0.993						

sample size fixed, we see fairly large changes in the upper bound on the rate of convergence of the contraction mapping.

The second column of Table 3 increases the standard deviation of the product-and-market-specific demand shocks, $\xi_{j,t}$. When these shocks are more variable, products become more vertically differentiated. Over the range of values we investigate, increases in the standard deviation of the demand shocks increase the Lipschitz constant. The third column of Table 3 changes the number of markets. The number of markets has little impact on the Lipschitz constant. Finally, the fourth column of Table 3 increases the mean of the intercept, $E[\beta_i^0]$, which changes the value of the inside goods relative to the outside good. As the inside good share increases, the Lipschitz constant increases.

7.4 Monte Carlo: Varying the Lipschitz Constant

Having established that different parameter settings can change the Lipschitz constant of the contraction mapping, we now explore whether there is an implication for execution time. We compare performances as we vary the mean of the intercept, $E[\beta_i^0]$, from -2 to 4. As we saw in Table 3, increasing $E[\beta_i^0]$ makes the Lipschitz constant higher. For each scenario, we run 20 replications of the data. For each data replication, we estimate the GMM parameters using our two numerically-accurate algorithms, NFP with a tight inner loop and MPEC. Because a local optimization routine may only converge to a local minimum, we follow what a rigorous researcher should do and we use multiple starting values for each algorithm and fake dataset. We run each algorithm five times per replication, using five independently-drawn starting values. We take the final point estimates for each algorithm as the run with the lowest objective function value. In all cases, the lowest objective function corresponded to a case where the algorithm reported that an optimal local solution had been found. We assess the estimates by looking at the own price elasticities, computed as a mean across

products within each market and then across markets. For each algorithm, we report the total CPU time required for all 10 runs. The results are reported in Table 4. All numbers in Table 4 are means across the 20 replications.

Turning to Table 4, we can see that our numerical theory prediction holds in practice. As expected, NFP with a tight inner loop tolerance and MPEC converge in all scenarios. We also find that MPEC and NFP generate identical point estimates, as one would expect since they are statistically the same estimator (Theorem 6). We compute the root mean-squared error (RMSE) and the bias of the own-price elasticities. For a parameter θ_1 , the bias is $E[\hat{\theta}_1] - \theta_1$, where θ_1 is the true value and the expectation is taking over many estimates with independent samples. Likewise, the RMSE is $\sqrt{E\left[\left(E[\hat{\theta}_1] - \theta_1\right)^2\right]}$. In all cases, both the bias and the RMSE are low, suggesting that the BLP estimator is capable of recovering true demand elasticities. To our knowledge, this is the most comprehensive Monte Carlo performed on BLP in the literature.

Run times vary dramatically for NFP tight with the level of the Lipschitz constant. For the low Lipschitz case with $E[\beta_i^0] = -2$, the average run time across the 20 replications (using 10 starting values for each replication) is roughly 20 minutes for NFP and 24 minutes for MPEC. However, as we increase the intercept, we see the run times for NFP increase, while the run times for MPEC change little. When $E[\beta_i^0] = 4$, the highest Lipschitz case, a single run of NFP takes, on average, 2.6 hours, whereas MPEC takes only 31 minutes. Thus, holding the number of inversions ($T \cdot J$) fixed, changing the intercept and, hence, the Lipschitz constant increased the CPU time for MPEC by a factor of 0.3, on average, and increased the CPU time for NFP by a factor of roughly 7 – several times more than for MPEC.

Thus, as predicted by the numerical theory, it is easy to find cases where NFP tight could be extremely slow to run due to the slow rate of convergence of the inner loop. In contrast, MPEC is fairly robust in terms of run times across scenarios. This relationship to run time highlights our earlier concern about the choice of inner loop tolerance. For real applications with many more products and/or markets (e.g. 25 products and 450 market/quarters in Nevo (2000, 2002) and 250 products and 10 years in BLP (1995)), run times could be considerably slower than in our Monte Carlo experiments with only 25 products and 50 markets. As we demonstrated previously, loosening the inner loop to speed the convergence of the inner loop could prevent the outer loop optimization from converging. This, in turn, might lead the researcher to loosen the outer loop tolerance, which could produce highly variable point estimates that may not even constitute local minima.

Our findings show the benefits of our suggested alternative estimation algorithm, MPEC. The MPEC algorithm offers several advantages. First, it eliminates the inner loop optimization and, hence, it is invariant to the rate of convergence of the contraction mapping. Furthermore, by eliminating

Table 4: Monte Carlo Results Varying the Lipschitz Constant

Intercept	Lipschitz	Implementation	Runs Converged	CPU Time (s)	Elasticities	
$E[\beta_i^0]$	Constant		(fraction)		Bias	RMSE
-2	0.806	NFP tight	1	1184.1	0.026	0.254
		MPEC	1	1455.1	0.026	0.254
-1	0.895	NFP tight	1	1252.8	0.029	0.258
		MPEC	1	1528.4	0.029	0.258
0	0.950	NFP tight	1	1352.5	0.029	0.265
		MPEC	1	1564.1	0.029	0.265
1	0.978	NFP tight	1	1641.1	0.029	0.270
		MPEC	1	1562.5	0.029	0.270
2	0.991	NFP tight	1	2498.1	0.030	0.273
		MPEC	1	1480.7	0.030	0.273
3	0.997	NFP tight	1	5128.1	0.031	0.276
		MPEC	1	1653.9	0.030	0.278
4	0.999	NFP tight	1	9248.5	0.032	0.279
		MPEC	1	1881.8	0.031	0.279

There are 20 replications for each experiment. Each replication uses five starting values to ensure a global minimum is found. The NFP tight implementation has $\epsilon_{\text{in}} = 10^{-14}$ and $\epsilon_{\text{out}} = 10^{-6}$. There is no inner loop in MPEC; $\epsilon_{\text{out}} = 10^{-6}$ and $\epsilon_{\text{feasible}} = 10^{-6}$. The same 100 simulation draws are used to generate the data and to estimate the model.

the inner loop, it avoids all the potential risks of naive implementations with loose tolerances. We therefore recommend MPEC as a safer and more reliable algorithm for the estimation of the BLP GMM estimator.

7.5 Varying the Number of Markets

In the previous section, we demonstrated that MPEC has a speed advantage over NFP when the Lipschitz constant is high. However, some readers may be concerned that MPEC may not be practical as one increases the number of products or the number of markets. The reason is that there is one nuisance optimization parameter, $\xi_{j,t}$, for each product j and market t combination. As the number of markets T (or the number of products J) increases, there will be more $\xi_{j,t}$ s over which to optimize and, correspondingly, more constraints. The next set of Monte Carlo experiments compare estimation with differing numbers of markets to see whether MPEC's speed advantage is related to having a small number of demand shocks.

Table 5 returns to the base specification, and varies only the number of markets, T . As the number of markets increases, not surprisingly both methods take longer. MPEC and NFP with tight tolerances take about the same amount of time until $T = 200$, at which point MPEC becomes faster. Although not reported in the table, at $T = 200$ MPEC is also faster than the naive implementations of NFP with loose inner-loop tolerances. For $T = 200$, MPEC takes roughly 42 minutes, whereas NFP with tight tolerances takes roughly 108 minutes. In contrast, for the base case with $T = 50$,

Table 5: Monte Carlo Results Varying the Number of Markets

# Markets	Lipschitz	Implementation	Runs Converged	CPU Time (s)
T	Constant		(fraction)	
25	0.937	NFP Tight	1	258.5
		MPEC	1	226.8
50 (base case)	0.944	NFP Tight	1	780.0
		MPEC	1	564.7
100	0.951	NFP Tight	1	2559.6
		MPEC	1	2866.0
200	0.953	NFP Tight	1	6481.7
		MPEC	1	2543.6

There are 20 replications for each experiment. Each replication uses five starting values to ensure a global minimum is found. The NFP tight implementation has $\epsilon_{\text{in}} = 10^{-14}$ and $\epsilon_{\text{out}} = 10^{-6}$. There is no inner loop in MPEC; $\epsilon_{\text{out}} = 10^{-6}$ and $\epsilon_{\text{feasible}} = 10^{-6}$. The same 100 simulation draws are used to generate the data and to estimate the model.

MPEC required only 9 minutes of CPU time and NFP required 13 minutes. Thus, increasing the number of markets increased MPEC’s CPU time by a factor of roughly 4.5, on average, whereas it increased NFP’s CPU time by a factor of roughly 8.3, on average – nearly double that of MPEC. We conclude that the performance advantages of MPEC over NFP actually increase as the number of demand shocks increase. Again, this result is not surprising. The modified Newton method used for MPEC has a quadratic rate of convergence whereas NFP has a linear rate of convergence for the inner loop. This means that MPEC should have a fairly easy time accommodating more parameters, in contrast with NFP accommodating a higher Lipschitz constant.

7.6 Speed Comparisons of MPEC and NFP Using Nevo’s Cereal Data

One potential criticism of our analysis above is that our Monte Carlo experiments were based on better-quality data than typical field data sets. In section (5.2), we used NFP with a tight tolerance to establish the reliability of the BLP GMM estimator for the cereal data. We now compare the speed of NFP and MPEC on this data. Like NFP, MPEC converged to the same local minimum with an objective function value of 4.5615 for 48 out of 50 starting values. For only two of the runs, MPEC converged to a different local minimum with a higher objective function value. In terms of run time for one starting value, we find that MPEC required an average CPU time of only 544 seconds whereas NFP required an average CPU time of 763 seconds. In short, the relative performance of MPEC and NFP documented in our Monte Carlo experiments appears to hold in the context of field data.

8 Other Computational Issues with BLP

8.1 Simulating Market Shares

The times for all methods reported in Table 4, the Monte Carlo results, are lower bounds on the actual speeds of these methods in applications. By shutting down simulation error, we were able to get by with $ns = 100$ simulation draws in the market share equations, (3). Our experiments with data generated using many more draws suggests that perhaps 10,000 draws might be appropriate to eliminate most simulation error, for models with five independent normal random coefficients. Using 10,000 instead of 100 draws will increase the speed of all three implementations by approximately $10,000/100 = 100$ times. The NFP algorithm with a tight tolerance for $E[\beta_i^0] = 4.0$, which in Table 4 took 9,248 seconds, would now take 924,800 seconds, or about 11 days. Because our main result regarding the speed advantage of MPEC would remain unchanged, we do not explore the role of simulation error in our Monte Carlo experiments.

8.2 Standard Errors

After obtaining point estimates, researchers need to compute standard errors to assess precision. Berry, Linton and Pakes (2004, Theorem 2) describe the sampling distribution of the BLP GMM estimator for fixed T and $J \rightarrow \infty$. As NFP and BLP are two computational implementations for the same estimator, both methods have the same sampling distribution.

One of the components of this formula requires derivatives of the mean of the moment conditions with respect to θ , or

$$\frac{\partial \left(\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^J \xi_{j,t}(\theta)' z_{j,t} \right)}{\partial \theta}.$$

This requires differentiating the inner loop. Berry, Linton, and Pakes (page 632) suggest numerical derivatives (finite derivatives) as one method of computing an estimate of this derivative. However, any attempt to numerically differentiate an inner loop has the potential to introduce substantial numerical error: numerical derivatives are often numerically inaccurate even when the function being differentiated itself has little numerical error. A more numerically accurate approach is to program the derivatives. These derivatives are found, for example, in the appendix of Nevo (2000b). The components of the NFP derivatives are also components of the MPEC derivatives, so coding the MPEC derivatives makes it easy to code the standard errors.²⁰ In the interests of simplicity and common-

²⁰For users who adopt MPEC, it may be possible to use the constrained distributions derived for GMM in Andrews (2002). Such a procedure requires simulating the asymptotic distribution: a realization of a normal random variable is drawn and then a constrained optimization problem is solved for each draw. We examined the results on extremum estimators with equality constraints in Gourieroux and Monfort (1995, Chapter 10). They derive the distribution of the constrained estimators as a function of the unconstrained estimators. For MPEC, the finite-sample objective function without constraints will always be minimized at $\xi = 0$. Therefore, the unconstrained limiting distribution is degenerate,

ality across researchers, we recommend users conducting $J \rightarrow \infty$ asymptotics code the asymptotic distribution in Theorem 2 of Berry, Linton and Pakes, using closed form derivatives. Users conducting $T \rightarrow \infty$ asymptotics and using a large enough number of simulation draws can use the standard GMM asymptotic variance formula.

8.3 Nonnegativity Constraints on Parameters

BLP (1995) and most subsequent empirical work uses a set of independent normal distributions for $F_\beta(\beta; \theta)$, the distribution of the random coefficients. Under normality, θ includes the standard deviation of each product characteristic’s random coefficient. The normal is a symmetric distribution. Therefore, if a guess for the standard deviation of characteristic 1’s random coefficient is σ_1 , $-\sigma_1$ should produce the same objective function value for NFP and both the same objective function and constraint values for MPEC. Any failure of this equivalence of σ_1 and $-\sigma_1$ under normality results from simulation error: (3) is not an accurate approximation to (2). Disregarding simulation error, the model is not identified unless the researcher constrains each standard deviation parameter to be nonnegative. If one of the standard deviation parameters is in truth zero, then Andrews (2002) shows how to conduct asymptotically valid hypothesis tests. The limiting distribution of the parameter on the boundary will be half-normal, as we know a standard deviation cannot be negative.

9 Extension: Maximum Likelihood Estimation

In this section, we outline how a researcher would adapt MPEC to a likelihood-based estimation of random coefficients logit demand. Some researchers prefer to work with likelihood-based estimators (Villas-Boas and Zhao 2005) and, more specifically, with Bayesian MCMC estimators (Jiang et al. 2008) based on the joint density of observed prices and market shares.²¹ Besides efficiency advantages, the ability to evaluate the likelihood of the data could be useful for testing purposes. The trade-off relative to GMM is the need for additional modeling structure which, if incorrect, could lead to biased parameter estimates. Like GMM, the calculation of the density of market shares still requires inverting the system of market share equations. Once again, MPEC can be used to circumvent the need for inverting the shares, thereby offsetting a layer of computational complexity and a potential source of numerical error. Below we outline the estimation of a limited information approach that models the data-generating process for prices in reduced form (much like two-stage least squares). However, one

and the proof technique in Gourieroux and Monfort does not apply.

²¹One can also think of Jiang et al. (2008) as an alternative algorithm for finding the parameters. The MCMC approach is a stochastic search algorithm that might perform well if the BLP model produces many local optima because MCMC will not be as likely to get stuck on a local flat region. Because our goal is not to study the role of multiple local minima, we do not explore the properties of a Bayesian MCMC algorithm.

can easily adapt the estimator to accommodate a structural (full-information) approach that models the data-generating process for supply-side variables, namely prices, as the outcome of an equilibrium in a game of imperfect competition (assuming the equilibrium exists and is unique).

Recall that the system of market shares is defined as follows:

$$s_j(x_t, p_t, \xi_t; \theta) = \int_{\beta} \frac{\exp(\beta^0 + x'_{j,t}\beta^x - \beta^p p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^0 + x'_{k,t}\beta^x - \beta^p p_{k,t} + \xi_{k,t})} dF_{\beta}(\beta; \theta). \quad (11)$$

We assume, as in a triangular system, that the data-generating process for prices is

$$p_{j,t} = z'_{j,t}\gamma + \eta_{j,t}, \quad (12)$$

where $z_{j,t}$ is a vector of price-shifting variables and $\eta_{j,t}$ is a mean-zero, i.i.d. shock. To capture the potential endogeneity in prices, we assume the supply and demand shocks have the following joint distribution: $(\xi_{j,t}, \eta_{j,t})' \equiv u_{j,t} \sim N(0, \Omega)$ where $\Omega = \begin{bmatrix} \sigma_{\xi}^2 & \sigma_{\xi, \eta} \\ \sigma_{\xi, \eta} & \sigma_{\eta}^2 \end{bmatrix}$.

The system defined by equations (11) and (12) has the joint density function

$$f_{s,p}(s_t, p_t; \Theta) = f_{\xi|\eta}(s_t | x_t, p_t; \theta, \Omega) |J_{\xi \rightarrow s}| f_{\eta}(p_t | z_t; \gamma, \Omega),$$

where $\Theta = (\theta, \gamma, \sigma_{\xi}^2, \sigma_{\xi, \eta}, \sigma_{\eta}^2)$ is the vector of model parameters, $f_{xi|\eta}(\cdot|\cdot)$ is the marginal density of ξ conditional on η , $f_{\eta}(\cdot|\cdot)$ is a Gaussian density with variance σ_{η}^2 , and $J_{\xi \rightarrow s}$ is the Jacobian matrix corresponding to the transformation of variables of $\xi_{j,t}$ to shares. The density of $\xi_{j,t}$ conditional on $\eta_{j,t}$ is

$$f_{\xi|\eta}(s_t | x_t, p_t; \theta, \Omega) = \prod_{j=1}^J \frac{1}{\sqrt{2\pi}\sigma_{\xi}\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2} \frac{(\xi_{j,t} - \rho \frac{\sigma_{\xi}}{\sigma_{\eta}} \eta_{j,t})^2}{\sigma_{\xi}^2 (1-\rho^2)}\right).$$

Note that the evaluation of $\xi_{j,t}$ requires inverting the market share equations, ((2)).

The element $J_{j,k}$ in row l and column k of the Jacobian matrix, $J_{\xi \rightarrow s}$, is

$$J_{j,l} = \begin{cases} \int_{\alpha, \beta} \left(1 - \frac{\exp(\beta^0 + x'_{j,t}\beta^x - \beta^p p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^0 + x'_{k,t}\beta^x - \beta^p p_{k,t} + \xi_{k,t})}\right) \frac{\exp(\beta^0 + x'_{j,t}\beta^x - \beta^p p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^0 + x'_{k,t}\beta^x - \beta^p p_{k,t} + \xi_{k,t})} dF(\theta) & , \quad j = l \\ - \int_{\alpha, \beta} \frac{\exp(\beta^0 + x'_{j,t}\beta^x - \beta^p p_{j,t} + \xi_{j,t})}{1 + \sum_{k=1}^J \exp(\beta^0 + x'_{k,t}\beta^x - \beta^p p_{k,t} + \xi_{k,t})} \frac{\exp(\beta^0 + x'_{l,t}\beta^x - \beta^p p_{l,t} + \xi_{l,t})}{1 + \sum_{k=1}^J \exp(\beta^0 + x'_{k,t}\beta^x - \beta^p p_{k,t} + \xi_{k,t})} dF(\theta) & , \quad j \neq l \end{cases}$$

Standard maximum likelihood estimation would involve searching for parameters, Θ^{LISML} , that

maximize the following log-likelihood function

$$l(\Theta) = \sum_{t=1}^T \log(f_{s,p}(s_t, p_t; \Theta)).$$

This would consist of a nested inner-loop to compute the demand shocks, $\xi_{j,t}$, via numerical inversion (the NFP contraction-mapping).

The equivalent MPEC approach entails searching for the vector of parameters (Θ, ξ) that maximizes the constrained optimization problem

$$\begin{aligned} l^{\text{MPEC}}(\Theta, \xi) &= \sum_{t=1}^T \log(f_{\xi|\eta}(s_t | x_t, p_t; \theta, \Omega) | J_{\xi \rightarrow s} | f_{\eta}(p_t | z_t; \gamma, \Omega)) \\ \text{subject to} & \quad s(\xi; \theta) = S \end{aligned} \tag{13}$$

10 Extension: Dynamic Demand Models

Starting with Melnikov (2000), a new stream of literature has considered dynamic analogs of BLP with forward-looking consumers making discrete choice purchases of durable goods (Nair 2007, Gordon 2007, Carranza 2008, Gowrisankaran and Rysman 2008, Dubé, Hitsch and Chintagunta 2008, Lee 2008, Schiraldi 2008). The typical implementation involves a nested fixed point approach with two nested inner loops. The first inner loop is the usual numerical inversion of the demand system to obtain the demand shocks, ξ . The second inner loop is the iteration of the Bellman equation to obtain the consumer's value function. In this section, we describe how MPEC can once again serve as a computationally more attractive solution than the nested fixed point approach.

As an example, we work with a simple model of demand for a durable good with falling prices over time. There is a mass M of potential consumers at date $t = 1$. Consumers are assumed to drop out of the market once they make a purchase. Abstracting from supply side specifics, we assume that prices evolve over time according to the rule

$$\log(p_{j,t}) = p'_{t-1}\rho_j + \psi_{j,t} \tag{14}$$

where $\psi_{j,t}$ is a random supply shock. For the remainder of our discussion, we assume that this supply shock is jointly-distributed with the demand shock: $(\xi_{j,t}, \psi_{j,t}) \sim N(0, \Omega)$ and is independent across time and markets. We assume that consumers have rational expectations in the sense that they use the true price process (14) to forecast future prices.

On the demand side, forward-looking consumers now have a real option associated with not purchasing because they can delay adoption to the future, when prices are expected to be lower. A

consumer r 's expected value of waiting is

$$\begin{aligned} v_0^r(p_t; \theta^r) &= \delta \int \max \left\{ \begin{array}{l} v_0^r(p'_t \rho_j + \psi; \theta^r) + \epsilon_0 \\ \max_j \{ \beta_j^r - \alpha^r (p'_t \rho_j + \psi + \psi) + \xi_j + \epsilon_j \} \end{array} \right\} dF_\epsilon(\epsilon) dF_{\psi, \xi}(\psi, \xi) \\ &= \delta \log \left(\exp(v_0^r(\mathcal{P}_j(p_t; \theta_p) + \psi; \theta^r)) + \sum_j \exp(\beta_j^r - \alpha^r (\mathcal{P}_j(p_t; \theta_p) + \psi) + \xi_j) \right) dF_{\psi, \xi}(\psi, \xi). \end{aligned} \quad (15)$$

To simplify the calculation of the expected value of waiting, we approximate it with Chebyshev polynomials (Judd 1998).²² We outline the Chebyshev approximation in Appendix C.

We use a discrete distribution to characterize the consumer population's tastes at date $t = 1$,

$$\theta^h \equiv \begin{bmatrix} \beta^h \\ \alpha^h \end{bmatrix} = \begin{cases} \theta^1, & \Pr(1) = \lambda_1 \\ \vdots & \vdots \\ \theta^R, & \Pr(R) = 1 - \sum_{r=1}^{R-1} \lambda_r \end{cases}.$$

This heterogeneity implies that certain types of consumers will systematically purchase earlier than others. Thus, the mass of remaining consumers of a given type r , M_t^r , evolves over time as follows:

$$M_t^r = \begin{cases} M \lambda_r & , t = 0 \\ M_{t-1}^r S_0^r(X_{t-1}; \Theta^r) & , t > 0 \end{cases}.$$

In a given period t , the market share of product j is

$$s_j(p_t; \theta) = \frac{\sum_{r=1}^R \lambda_{t,r} \frac{\exp(\beta_j^r - \alpha^r p_{j,t} + \xi_{j,t})}{\exp(v_0^r(p_t; \theta^r)) + \sum_{k=1}^J \exp(\beta_k^r - \alpha^r p_{k,t} + \xi_{k,t})}}{\sum_{r=1}^R \lambda_{t,r} \frac{\exp(\beta_j^r - \alpha^r p_{j,t} + \xi_{j,t})}{\exp(v_0^r(p_t; \theta^r)) + \sum_{k=1}^J \exp(\beta_k^r - \alpha^r p_{k,t} + \xi_{k,t})}}, \quad (16)$$

where

$$\lambda_{t,r} = \begin{cases} \lambda_r & t = 0 \\ \frac{M_t^r}{\sum_r M_t^r} & , t > 0 \end{cases}$$

is the proportion of type r consumers still in the market at date t .

The empirical model consists of the system (14) and (16), which we write more compactly as

$$u_t \equiv \begin{bmatrix} \psi_t \\ \xi_t \end{bmatrix} = \begin{bmatrix} \log(p_t) - p'_{t-1} \rho \\ s^{-1}(p_t, S_t; \Theta) \end{bmatrix}.$$

We use the joint density of $(\xi_{j,t}, \psi_{j,t})$ to construct a maximum likelihood estimator of the model

²²One could also add 0.577 to the expected value of waiting in order to account for the mean of the logit error term.

parameters. The multivariate normal distribution of $(\xi_{j,t}, \psi_{j,t})$ induces the density on the observable outcomes, (p, S_t) ,

$$f_{p,S}(p_t, S_t; \theta, \rho, \Omega) = \frac{1}{(2\pi)^{\frac{3J}{2}} |\Omega|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} u_t' \Omega_u^{-1} u_t\right) |J_{t,u \rightarrow Y}|$$

where $J_{t,u \rightarrow Y}$ is the $(2J \times 2J)$ Jacobian matrix corresponding to the transformation-of-variables from u_t to Y_t . We provide the derivation of the Jacobian in Appendix D.

An NFP approach to estimating the model parameters amounts to solving the optimization problem

$$\max_{\{\theta, \rho, \Omega\}} \prod_{t=1}^T f_{p,S}(p_t, S_t; \theta, \rho, \Omega). \quad (17)$$

This problem nests two inner loops. For each stage of the outer loop to maximize the likelihood function in (17), one needs to solve for a fixed point of the contraction mapping, (15), in order to obtain the expected value of waiting. In addition, one needs to solve the fixed point of the BLP contraction mapping, 5, to compute the demand shocks ξ_t (i.e. the inversion). Numerical error from both these inner loops can potentially propagate into the outer loop. Thus, the numerical concerns regarding inner loop convergence tolerance discussed for static BLP are exacerbated with dynamic analogs of BLP.

Let D be the support of the state variables. An MPEC approach to estimating the model parameters amounts to solving the optimization problem

$$\begin{aligned} \max_{\{\theta, \rho, \Omega, \xi, v\}} & \prod_{t=1}^T \frac{1}{(2\pi)^{\frac{3J}{2}} |\Omega|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} u_t' \Omega_u^{-1} u_t\right) |J_{t,u \rightarrow Y}| \\ \text{subject to} & \quad s(\xi_t; \theta) = S_t \quad \forall t = 1, \dots, T \\ \text{and} & \quad v_0^r(p_d) = \delta \log \left(\frac{\exp(v_0^r(p_d' \rho_j + \psi)) + \dots}{\sum_j \exp(\beta_j^r - \alpha^r (p_d' \rho_j + \psi) + \xi_j)} \right) dF_{\psi, \xi}(\psi, \xi) \\ & \quad \forall d \in D, r = 1, \dots, R. \end{aligned}$$

In this formulation, we now optimize over the demand shocks, ξ , and the expected value of waiting evaluated at each point, $v^r(p_d)$. In this case $D \subset \mathbb{R}_+^2$, which is the support of the two products' prices. While this approach increases the number of parameters in the outer-loop optimization problem substantially compared to NFP, MPEC completely eliminates the two inner loops. Note that the use of Chebyshev approximation reduces the dimension of this problem substantially. Rather than searching over the value function at each point in a discretized state space, we search over the Chebyshev weights.

To assess the relative performance of MPEC versus NFP in the context of our dynamic durable goods example, we construct the following Monte Carlo experiments. In the first experiment, we

assume there is only a single consumer type, $R = 1$. It is easy to show that in this case, ξ_t can be computed analytically by log-linearizing the market shares, (16).²³ We begin with this case because it only involves a nested call to the calculation of the expected value of waiting. Below we will allow for more consumer types to see what happens when we also require a nested call to the numerical inversion of the shares. We assume that the consumers' preferences are: $(\beta_1, \beta_2, \alpha) = (4, -1, -.15)$ and the discount factor is $\delta = 0.99$.²⁴ We assume that the density of prices has the transition rules

$$\begin{bmatrix} p_{1,t} & = & 5 + .8p_{1,t-1} + .2p_{2,t-1} + \psi_{1,t} \\ p_{2,t} & = & 5 + .1p_{1,t-1} + .55p_{2,t-1} + \psi_{2,t} \end{bmatrix}.$$

Note how the lagged price of product 2 effects the price of product 1, and vice versa. Finally, we assume the supply and demand shocks satisfy $(\psi_{j,t}, \xi_{j,t}) \sim N\left(0, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right)$ and are independent across markets and time periods. For our Chebyshev approximation, we use 6 grid points and a 6th order polynomial. For the NFP algorithm, we use an inner loop tolerance of 10^{-14} for the calculation of the expected value of waiting.

Results from 25 replications of this first experiment are reported in Table 6. We report the bias and RMSE associated with each of the structural parameters, for MPEC and NFP respectively. Interestingly, MPEC seems to produce estimates that, on average, have lower bias while NFP seems to produce lower RMSE. This may be a consequence of using only one starting value per replication. More importantly, the average CPU time for MPEC is just over 25% of the CPU time for NFP.

Now we run a second Monte Carlo experiment where we allow for two types of consumers.

11 Conclusions

In this paper, we analyzed the numerical properties of the NFP approach proposed by BLP to estimate the random coefficients logit demand model. Theoretically, the NFP approach may be slow, as NFP's inner loop is only linearly convergent and NFP is more vulnerable to error due to the inner loop. We showed the Lipschitz constant is a measure of an upper bound to the convergence rate of NFP's inner loop's contraction mapping. We numerically evaluated the Lipschitz constant for particular data generating processes and showed when the inner loop is likely to be slow. Further, we showed that setting loose inner loop tolerances can lead to incorrect parameter estimates and a failure of the optimization routine to report that it has converged.

We then proposed a new constrained optimization formulation, MPEC, for estimating the random

²³See Berry (1994) on how to invert the demand shocks in the homogeneous logit model.

²⁴A discount factor of 0.99 at the quarterly level corresponds to an annual discount rate of 0.96, which is a commonly-used value in the literature.

	Bias		RMSE	
	MPEC	NFP	MPEC	NFP
θ : taste parameters				
$\beta_1 : 4$	7.5E-03	4.6E-02	1.7E-01	1.5E-01
$\beta_2 : -1$	6.2E-03	3.7E-02	1.5E-01	1.2E-01
$\alpha : -0.15$	-1.1E-04	-2.9E-04	8.0E-04	5.4E-04
ρ : price transitions				
$int_1 : 5$	9.4E-03	1.9E-02	4.9E-02	4.6E-02
$\rho_{1,1} : 0.8$	9.5E-05	-2.1E-04	1.2E-03	1.2E-03
$\rho_{1,2} : 0.2$	-1.6E-04	-3.8E-05	1.5E-03	1.7E-03
$int_2 : 5$	8.9E-03	6.6E-04	5.9E-02	3.2E-02
$\rho_{2,1} : 0.1$	-7.0E-05	1.5E-04	1.1E-03	5.6E-04
$\rho_{2,2} : 0.55$	-6.5E-05	-4.5E-04	1.4E-03	8.8E-04
Ω : variances of shocks				
1	-4.1E-03	-4.5E-03	1.7E-02	1.7E-02
0.866	-1.7E-03	-5.5E-04	1.5E-02	1.4E-02
0.5	-7.9E-04	-2.4E-03	2.0E-02	1.9E-02
Avg CPU time (sec)	4579.3	16,971		

Table 6: Monte Carlo Results for Dynamic BLP with One Consumer Type: NFP versus MPEC

coefficient logit demand model. MPEC is quicker to compute and avoids numerical errors because it avoids repeatedly inverting the market shares equations numerically. It also allows the researcher to access state-of-the-art constrained optimization solvers.

To assess the practical aspects of MPEC versus NFP, we conducted a number of Monte Carlo experiments. In many instances, we find that the NFP approach can produce accurate estimates of the demand parameters in a reasonable amount of time. However, we find that NFP is comparatively quite slow under a few data generating processes. The Lipschitz constant is a measure of an upper bound to the NFP inner loop’s speed. Applications with high Lipschitz constants might take a long time for estimation. Loosening the inner loop tolerance to improve the speed of NFP, as is often done in practice, leads to a failure to converge unless the optimization method’s tolerance is also significantly loosened, which can potentially be dangerous.

In contrast, MPEC produces good estimates relatively quickly for all data generating processes. The reason is that there is no inner loop in MPEC, so issues like a high Lipschitz constant (or slow contraction mapping) are irrelevant for MPEC. MPEC takes about the same amount of time for data generating processes with different Lipschitz constants and hence different NFP speeds. MPEC is also useful for estimating the BLP demand model using maximum likelihood, where a Jacobian term involving the demand shocks must be computed.

As an extension, we adapt the MPEC approach to a new class of applications with forward-looking consumers. The relative advantage of MPEC is even stronger with dynamics because two inner loops must be solved: the dynamic programming problem and the market share inversion. This burdensome

collection of three loops (optimization, market shares, dynamic programming) makes the traditional BLP approach nearly untenable in terms of computational time. Current work (Lee 2008, Schiraldi 2008) further extends the number of inner loops being solved in estimation. As demand models become richer, the computational benefits of MPEC over NFP become greater.

References

- [1] Akerberg, D. and M. Rysman (2005), Unobserved Product Differentiation in Discrete Choice Models: Estimating Price Elasticities and Welfare Effects, *RAND Journal of Economics*, 36 (4): 771–788.
- [2] Andrews, D. W. K. (2002). Generalized Method of Moments Estimation When a Parameter is on a Boundary, *Journal of Business and Economic Statistics*, 20 (4), 530–544.
- [3] Bajari, P., J. T. Fox and S. P. Ryan (2007), “Linear Regression Estimation of Discrete Choice Models with Nonparametric Distributions of Random Coefficients.” *American Economic Review*, 97(2), 459-463
- [4] Bajari, P., J. T. Fox, K.-I. Kim and S. P. Ryan (2008), A Simple Nonparametric Estimator for the Distribution of Random Coefficients in Discrete Choice Models. Working paper.
- [5] Berry, S. (1994), Estimating Discrete-Choice Models of Product Differentiation. *RAND Journal of Economics*, 25(2): 242–262.
- [6] Berry, S. and P. A. Haile (2008), Nonparametric Identification of Multinomial Choice Models with Heterogeneous Consumers and Endogeneity, working paper.
- [7] Berry, S., J. Levinsohn, and A. Pakes (1995), Automobile Prices in Market Equilibrium. *Econometrica*, 63(4): 841–890.
- [8] Berry, S., J. Levinsohn, and A. Pakes (2004), “Differentiated Products Demand Systems from a Combination of Micro and Macro Data: The New Car Market.” *Journal of Political Economy*, 112(1):68-105.
- [9] Berry, S., O. B. Linton, . and A. Pakes (2004), Limit Theorems for Estimating the Parameters of Differentiated Product Demand Systems. *Review of Economic Studies*, 71(3): 613–654.
- [10] Byrd, R. H., J. Nocedal and R. A. Waltz (1999), An Interior Point Method for Large Scale Nonlinear Programming. *SIAM J. Optimization*, 9, 4, 877-990.

- [11] Carranza, Juan Esteban. (2008) Product innovation and adoption in market equilibrium: The case of digital cameras. Working paper.
- [12] Dahlquist, Germund and Åke Björck. (2008) *Numerical Methods in Scientific Computing*. SIAM, Philadelphia.
- [13] Davis, Peter J. (2006) "The Discrete Choice Analytically Flexible (DCAF) Model of Demand for Differentiated Products," CEPR Discussion Papers 5880.
- [14] Dubé, Jean-Pierre, Guenter Hitsch and Pradeep Chintagunta (2008). "Tipping and Concentration in Markets With Indirect Network Effects," Working paper, University of Chicago.
- [15] Fox, Jeremy T. and Amit Gandhi (2008) Identifying Heterogeneity in Economic Choice and Selection Models Using Mixtures. Working paper.
- [16] Gandhi, Amit (2008). "On the Nonparametric Foundations of Product Differentiated Demand Systems". Working paper.
- [17] Gourieroux, Christian and Alain Monfort (1995), *Statistics and Econometric Models*, Vol 1. Cambridge.
- [18] Gowrisankaran, G. and M. Rysman (2007), Dynamics of Consumer Demand for New Durable Goods. Working paper.
- [19] Hausman, J.A. and D. A. Wise (1976). "A Conditional Profit Model for Qualitative Choice: Discrete Decisions Recognizing Interdependence and Heterogeneous Preferences." *Econometrica*, 46(2), 403–426.
- [20] Hendel, I. and Aviv Nevo (2007) Measuring the Implications of Sales and Consumer Inventory Behavior, *Econometrica*, 74(16), 1637-1673.
- [21] Jiang, Renna, Puneet Manchanda and Peter E. Rossi (2007), "Bayesian Analysis of Random Coefficient Logit Models Using Aggregate Data." Working Paper.
- [22] Judd, Kenneth L. (1992) "Projection methods for solving aggregate growth models." *Journal of Economic Theory*, 58(2), 410–452.
- [23] K. L. Judd (1998), *Numerical Methods in Economics*. MIT Press.
- [24] C. T. Kelley (1995), *Iterative Methods for Linear and Nonlinear Equations*. SIAM, Philadelphia.
- [25] C. T. Kelley (1999), *Iterative Methods for Optimization*, SIAM, Philadelphia.

- [26] C. T. Kelley (2003), *Solving Nonlinear Equations with Newton's Method*. SIAM, Philadelphia.
- [27] Knittel, Christopher R. and Konstantinos Metaxoglou (2008), "Estimation of Random Coefficient Demand Models: Challenges, Diculties and Warnings." Working paper.
- [28] Lee, Robin S. "Vertical Integration and Exclusivity in Platform and Two-Sided Markets." Working paper.
- [29] Mas-Colell, Andreu, Michael D. Whinston, and Jerry R. Green. (1995) *Microeconomic Theory*. Oxford University Press.
- [30] McFadden, D. and K. Train (2000), Mixed MNL models for discrete response. *Journal of Applied Econometrics*, 15(5): 447–470.
- [31] Melnikov, Oleg (2001) Demand for Differentiated Durable Products: The Case of the U.S. Computer Printer Market. Working paper.
- [32] Nair, Harikesh. Intertemporal Price Discrimination with Forward-looking Consumers: Application to the US Market for Console Video-Games. *Quantitative Marketing and Economics*, 5(3), 239-292.
- [33] Nevo, A. (2000a), Mergers with Differentiated Products: The Case of the Ready-to-Eat Cereal Industry. *RAND Journal of Economics*, 31(3): 395–421.
- [34] Nevo, A.(2000b), A Practitioner's Guide to Estimation of Random Coefficients Logit Models of Demand, *Journal of Economics and Management Strategy*, 9(4): 513–548.
- [35] Nevo, A. (2001), Measuring Market Power in the Ready-to-Eat Cereal Industry. *Econometrica*, 69(2): 307–342.
- [36] Nocedal, J. and S. J. Wright (2006), *Numerical Optimization*. Springer, New York.
- [37] Petrin, A. (2002), Quantifying the Benefits of New Products: The Case of the Minivan. *Journal of Political Economy*, 110:705–729.
- [38] Petrin, A. and K. Train (2006), Control Function Corrections for Omitted Attributes in Differentiated Products Markets. Working paper.]
- [39] Schiraldi, Pasquale (2008). "Automobile Replacement: a Dynamic Structural Approach." Working paper.
- [40] Su, C.-L. and K. L. Judd (2007), Constrained Optimization Approaches to Estimation of Structural models. Working paper, CMS-EMS, Kellogg School of Management.

- [41] van der Vaart, A. W. (2000). *Asymptotic Statistics*. Cambridge.
- [42] Villas-Boas, J. Miguel and Russell S. Winer (1999). "Endogeneity in Brand Choice Models." *Management Science*, 45:1324-1338.
- [43] Villas-Boas, J. Miguel and Ying Zhao (2005). "Retailer, Manufacturers, and Individual Consumers: Modeling the Supply Side in the Ketchup Marketplace." *Journal of Marketing Research*, 42, 83-95.

A Proofs

A.1 Proof of Theorem 3

By a Taylor series expansion of $Q(\xi)$ around $\xi(\theta, 0)$, we have

$$\begin{aligned} & Q(\xi(\theta, \epsilon_{\text{in}})) - Q(\xi(\theta, 0)) \\ &= \left[\frac{\partial Q(\xi)}{\partial \xi} \Big|_{\xi=\xi(\theta, 0)} \right]' (\xi(\theta, \epsilon_{\text{in}}) - \xi(\theta, 0)) + O\left(\|\xi(\theta, \epsilon_{\text{in}}) - \xi(\theta, 0)\|^2\right) \\ & \text{and} \\ & \nabla_{\theta} Q(\xi) \Big|_{\xi=\xi(\theta, \epsilon_{\text{in}})} - \nabla_{\theta} Q(\xi) \Big|_{\xi=\xi(\theta, 0)} \\ &= \left[\frac{\partial \nabla_{\theta} Q(\xi(\theta))}{\partial \xi} \Big|_{\xi=\xi(\theta, 0)} \right]' (\xi(\theta, \epsilon_{\text{in}}) - \xi(\theta, 0)) + O\left(\|\xi(\theta, \epsilon_{\text{in}}) - \xi(\theta, 0)\|^2\right). \end{aligned}$$

Because $\|\xi(\theta, \epsilon_{\text{in}}) - \xi(\theta, 0)\| \leq \frac{L(\theta)}{1-L(\theta)} \epsilon_{\text{in}}$ by Theorem 1, and assuming both $\left\| \frac{\partial Q(\xi)}{\partial \xi} \Big|_{\xi=\xi(\theta, 0)} \right\|$ and $\left\| \frac{\partial \nabla_{\theta} Q(\xi(\theta))}{\partial \xi} \Big|_{\xi=\xi(\theta, 0)} \right\|$ are bounded, we obtain

$$\begin{aligned} |Q(\xi(\theta, \epsilon_{\text{in}})) - Q(\xi(\theta, 0))| &= O\left(\frac{L(\theta)}{1-L(\theta)} \epsilon_{\text{in}}\right) \\ \|\nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\theta, \epsilon_{\text{in}})} - \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\theta, 0)}\| &= O\left(\frac{L(\theta)}{1-L(\theta)} \epsilon_{\text{in}}\right). \end{aligned}$$

A.2 Proof of Theorem 4

We define θ^* to be the true estimate when there are no inner-loop numerical errors ($\epsilon_{\text{in}} = 0$) and $\hat{\theta}$ the numerically incorrect estimates with the inner-loop tolerance ϵ_{in} ,

$$\begin{aligned} \theta^* &= \arg \max_{\theta} \{Q(\xi(\theta, 0))\} \\ \hat{\theta} &= \arg \max_{\theta} \{Q(\xi(\theta, \epsilon_{\text{in}}))\}. \end{aligned}$$

Because $\nabla_{\theta} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, \epsilon_{\text{in}})} = 0$, the application of the second result in Theorem 3 at $\hat{\theta}$ gives

$$\left\| \nabla_{\theta} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right\| = O\left(\frac{L(\hat{\theta})}{1-L(\hat{\theta})} \epsilon_{\text{in}}\right).$$

Note that we have evaluated the GMM objective function with no numerical error at the point $\hat{\theta}$ that minimizes the GMM objective function with inner loop numerical error.

Let $\tilde{\theta}$ be any value of the structural parameters near $\hat{\theta}$. By first the inverse triangle inequality,

then the regular triangle inequality, and then finally a Taylor series expansion, we have

$$\begin{aligned}
& \left\| \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, \epsilon_{\text{in}})} \right\| - \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right\| \right\| \\
\leq & \left\| \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, \epsilon_{\text{in}})} \right\| - \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, 0)} \right\| \right\| \\
\leq & \left\| \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, \epsilon_{\text{in}})} - \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, 0)} \right\| \right\| \\
= & \left\| \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, \epsilon_{\text{in}})} - \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, 0)} + \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, 0)} - \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right\| \right\| \\
\leq & \left\| \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, \epsilon_{\text{in}})} - \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, 0)} \right\| \right\| \\
& + \left\| \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, 0)} - \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right\| \right\| \\
\leq & O\left(\frac{L(\hat{\theta})}{1-L(\hat{\theta})}\epsilon_{\text{in}}\right) + \left\| \nabla_{\theta}^2 Q(\xi(\theta)) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right\| \|\tilde{\theta} - \hat{\theta}\| + O\left(\|\tilde{\theta} - \hat{\theta}\|^2\right).
\end{aligned}$$

As we have assumed $\left\| \nabla_{\theta}^2 Q(\xi(\theta)) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right\| \|\tilde{\theta} - \hat{\theta}\|$ is bounded, we obtain

$$\begin{aligned}
\left\| \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, \epsilon_{\text{in}})} \right\| \right\| & \leq \left\| \left\| \nabla_{\theta} Q(\xi(\theta)) \Big|_{\xi=\xi(\tilde{\theta}, 0)} \right\| \right\| + O\left(\frac{L(\hat{\theta})}{1-L(\hat{\theta})}\epsilon_{\text{in}}\right) + O\left(\|\tilde{\theta} - \hat{\theta}\|^2\right) \\
& = O\left(\frac{L(\hat{\theta})}{1-L(\hat{\theta})}\epsilon_{\text{in}}\right) + O\left(\|\tilde{\theta} - \hat{\theta}\|^2\right).
\end{aligned}$$

Hence, the norm of the numerically inaccurate gradient, evaluated at an arbitrary point $\tilde{\theta}$, is bounded above by a term on the order of $\frac{L(\hat{\theta})}{1-L(\hat{\theta})}\epsilon_{\text{in}}$ and a term involving the arbitrary point $\tilde{\theta}$ and the GMM estimator with error, $\hat{\theta}$. The term $O\left(\frac{L(\hat{\theta})}{1-L(\hat{\theta})}\epsilon\right)$ indicates that the numerical error in the gradient is linearly increasing in ϵ_{in} (decreasing ϵ_{in} decreases the numerical error in the gradient).

A.3 Proof of Theorem 5

First, we can quantify the bias between the numerically correct and incorrect objective function values, $Q(\xi(\hat{\theta}, \epsilon_{\text{in}}))$ and $Q(\xi(\theta^*, 0))$. By two Taylor series expansions, we have

$$\begin{aligned}
& Q(\xi(\hat{\theta}, \epsilon_{\text{in}})) - Q(\xi(\theta^*, 0)) \\
= & Q(\xi(\hat{\theta}, \epsilon_{\text{in}})) - Q(\xi(\hat{\theta}, 0)) + Q(\xi(\hat{\theta}, 0)) - Q(\xi(\theta^*, 0)) \\
= & \left[\nabla_{\xi} Q(\xi(\theta)) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right]' (\xi(\hat{\theta}, \epsilon_{\text{in}}) - \xi(\hat{\theta}, 0)) + O\left(\|\xi(\hat{\theta}, \epsilon_{\text{in}}) - \xi(\hat{\theta}, 0)\|^2\right) + \\
& \left[(\nabla_{\theta} \xi(\theta))' \nabla_{\xi} Q(\xi) \Big|_{\xi=\xi(\theta^*, 0)} \right]' (\hat{\theta} - \theta^*) + O\left(\|\hat{\theta} - \theta^*\|^2\right) \\
= & \left[\nabla_{\xi} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, 0)} \right]' (\xi(\hat{\theta}, \epsilon_{\text{in}}) - \xi(\hat{\theta}, 0)) + O\left(\|\xi(\hat{\theta}, \epsilon_{\text{in}}) - \xi(\hat{\theta}, 0)\|^2\right) + O\left(\|\hat{\theta} - \theta^*\|^2\right),
\end{aligned}$$

because $\nabla \xi(\theta^*)' \nabla_{\xi} Q(\xi(\theta^*)) = 0$ at the true estimates θ^* .

Rearranging the equality involving $Q(\xi(\hat{\theta}, \epsilon_{in})) - Q(\xi(\theta^*, 0))$ to focus on the $O(\|\hat{\theta} - \theta^*\|^2)$ term, we have

$$\begin{aligned} O(\|\hat{\theta} - \theta^*\|^2) &= Q(\xi(\hat{\theta}, \epsilon_{in})) - Q(\xi(\theta^*, 0)) - [\nabla_{\xi} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, 0)}]' (\xi(\hat{\theta}, \epsilon_{in}) - \xi(\hat{\theta}, 0)) \\ &\quad - O(\|\xi(\hat{\theta}, \epsilon_{in}) - \xi(\hat{\theta}, 0)\|^2) \\ &\leq |Q(\xi(\hat{\theta}, \epsilon_{in})) - Q(\xi(\theta^*, 0))| + \|\nabla_{\xi} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, 0)}\| \|\xi(\hat{\theta}, \epsilon_{in}) - \xi(\hat{\theta}, 0)\| \\ &\quad - O(\|\xi(\hat{\theta}, \epsilon_{in}) - \xi(\hat{\theta}, 0)\|^2). \end{aligned}$$

To use this bound numerically, assume that $\|\nabla_{\xi} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, 0)}\|$ is bounded and that

$$O(\|\xi(\hat{\theta}, \epsilon_{in}) - \xi(\hat{\theta}, 0)\|^2) \ll |Q(\xi(\hat{\theta}, \epsilon_{in})) - Q(\xi(\theta^*, 0))| + \|\nabla_{\xi} Q(\xi) \Big|_{\xi=\xi(\hat{\theta}, 0)}\| \|\xi(\hat{\theta}, \epsilon_{in}) - \xi(\hat{\theta}, 0)\|.$$

This allows us to focus on the numerical error from the NFP algorithm's inner loop and the bias in objective values. We also know from the choice of the contraction mapping inner loop tolerance that $\|\xi(\hat{\theta}, \epsilon_{in}) - \xi(\hat{\theta}, 0)\| \leq \frac{L(\hat{\theta})}{1-L(\hat{\theta})} \epsilon_{in}$. Therefore, $\|\xi(\hat{\theta}, \epsilon_{in}) - \xi(\hat{\theta}, 0)\|$ is also bounded. Hence, we obtain

$$O(\|\hat{\theta} - \theta^*\|^2) \leq |Q(\xi(\hat{\theta}, \epsilon_{in})) - Q(\xi(\theta^*, 0))| + O\left(\frac{L(\hat{\theta})}{1-L(\hat{\theta})} \epsilon_{in}\right).$$

A.4 Proof of Theorem 6

The NFP method (4) solves the following unconstrained problem

$$\min_{\theta} Q(\xi(\theta)). \tag{18}$$

The first-order condition of (18) is

$$\frac{\partial Q(\xi(\theta))}{\partial \theta} = \frac{d\xi'}{d\theta} \frac{\partial Q}{\partial \xi} = 0. \tag{19}$$

The constrained optimization formulation of (18) is

$$\begin{aligned} \min_{(\theta, \xi)} \quad & Q(\xi) \\ \text{s.t.} \quad & s(\xi; \theta) = S. \end{aligned} \tag{20}$$

The Lagrangian for (20) is $\mathcal{L}(\theta, \xi, \lambda) = Q(\xi) - \lambda^T (S - s(\xi; \theta))$, where λ is the vector of Lagrange multipliers. The first-order conditions of (20) are

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta, \xi, \lambda)}{\partial \theta} &= - \frac{ds(\xi; \theta)'}{d\theta} \lambda = 0 \\ \frac{\partial \mathcal{L}(\theta, \xi, \lambda)}{\partial \xi} &= \frac{\partial Q}{\partial \xi} - \frac{ds(\xi; \theta)'}{d\xi} \lambda = 0 \\ \frac{\partial \mathcal{L}(\theta, \xi, \lambda)}{\partial \lambda} &= S - s(\xi; \theta) = 0.\end{aligned}\tag{21}$$

Since the BLP inner loop is a contraction mapping, the matrix $\frac{ds(\xi; \theta)'}{d\xi}$ is invertible.²⁵ Solving the second set of first order conditions for λ gives $\lambda = \left(\frac{ds(\xi; \theta)'}{d\xi}\right)^{-1} \frac{\partial Q}{\partial \xi}$. Then

$$\frac{\partial \mathcal{L}}{\partial \theta} = - \frac{ds(\xi; \theta)'}{d\theta} \left(\frac{ds(\xi; \theta)'}{d\xi}\right)^{-1} \frac{\partial Q}{\partial \xi} = 0,\tag{22}$$

which is identical to (19), the first-order condition from the NFP formulation. To see the equivalence, note that the implicit function theorem (Theorem M.E.1 in Mas-Colell, Whinston and Green 1995) states

$$\frac{\partial \xi(\theta)}{\partial \theta} = - \left(\frac{ds(\xi; \theta)'}{d\xi}\right)^{-1} \frac{ds(\xi; \theta)'}{d\theta},$$

so by substitution,

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \xi(\theta)'}{\partial \theta} \frac{\partial Q}{\partial \xi} = \frac{\partial Q(\xi(\theta))}{\partial \theta}.$$

B Gradients for the MPEC Objective Function and Constraints

Here we derive the gradients of the MPEC objective function and constraints with respect to the optimization parameters in MPEC. These gradients are an important input, for both numerical accuracy and speed. Nevo (2000) lists the gradients for NFP. This section uses the independent normal distribution for each of the random coefficients, as in BLP (1995) and many other empirical papers.

Market Share

$$\begin{aligned}s_j(\xi_t; \theta) &= \int \frac{\exp(x'_{j,t} \bar{\beta} - \bar{\alpha} p_{j,t} + \xi_{j,t} + \sum_k x'_{k,j,t} \nu_k \sigma_{\beta_k} - p_{j,t} \nu_{K+1} \sigma_{\alpha})}{1 + \sum_{i=1}^J \exp(x'_{i,t} \bar{\beta} - \bar{\alpha} p_{i,t} + \xi_{i,t} + \sum_k x'_{k,i,t} \nu_k \sigma_{\beta_k} - p_{i,t} \nu_{K+1} \sigma_{\alpha})} dF(\nu) \\ &= \int T_j(\xi_t, \nu; \theta) dF(\nu)\end{aligned}$$

where $\theta = (\bar{\beta}, \bar{\alpha}, \sigma_{\beta}, \sigma_{\alpha})'$, and $\nu \sim N(0, I_{K+1})$.

MPEC Criterion Function

²⁵We thank Ken Judd and John Birge for pointing out this property.

$$\begin{aligned}
& \min_{\theta, \xi} && g(\xi)' W g(\xi) \\
& \text{subject to} && s(\xi; \theta) = S
\end{aligned} \tag{23}$$

$$\text{where } g(\xi) = \frac{1}{T} \sum_{t=1}^T \xi_t' z_t$$

Gradients for MPEC

$$\frac{\partial s_j(\xi_t; \theta)}{\partial \beta_k} = \int T_j(\xi_t, \nu; \theta)(x_{j,k,t} - \sum_i T_i(\xi_t, \nu; \theta)x_{k,i,t})dF(\nu)$$

$$\frac{\partial s_j(\xi_t; \theta)}{\partial \bar{\alpha}} = \int T_j(\xi_t, \nu; \theta)(p_{j,k,t} - \sum_i T_i(\xi_t, \nu; \theta)p_{k,i,t})dF(\nu)$$

$$\frac{\partial s_j(\xi_t; \theta)}{\partial \sigma_{\beta_k}} = \int T_j(\xi_t, \nu; \theta)(x_{j,k,t} - \sum_i T_i(\xi_t, \nu; \theta)x_{k,i,t})\nu_k dF(\nu)$$

$$\frac{\partial s_j(\xi_t; \theta)}{\partial \sigma_{\alpha}} = \int T_j(\xi_t, \nu; \theta)(p_{j,k,t} - \sum_i T_i(\xi_t, \nu; \theta)p_{k,i,t})\nu_{K+1} dF(\nu)$$

$$\frac{\partial s_j(\xi_t; \theta)}{\partial \xi_{j,t}} = \int T_j(\xi_t, \nu; \theta)(1 - T_j(\xi_t, \nu; \theta))dF(\nu)$$

$$\frac{\partial s_j(\xi_t; \theta)}{\partial \xi_{i,t}} = - \int T_j(\xi_t, \nu; \theta)T_i(\xi_t, \nu; \theta)dF(\nu)$$

$$\frac{\partial g(\xi)' W g(\xi)}{\partial \xi} = 2g(\xi)' W \frac{\partial g(\xi)}{\partial \xi}$$

C Chebyshev Approximation of the Expected Value of Waiting

First, we bound the range of prices as follows, $p = (p_1, p_2)' \in [0, b] \times [0, b]$, where b is large (b is 1.5 times the largest observed price in the data). We then approximate the expected value of delaying adoption with Chebyshev polynomials, $v_0^r(p; \theta^r) \approx \gamma^{r'} \Lambda(p)$, where γ^r is a $K \times 1$ vector of parameters and $\Lambda(p)$ is a $K \times 1$ vector of K Chebyshev polynomials. Therefore, we can re-write the Bellman as

$$\gamma^{r'} \Lambda(p) = \delta \int \log \left(\exp \left(\gamma^{r'} \Lambda(p\rho + \psi) \right) + \sum_j \exp \left(\beta_j^r - \alpha^r (p' \rho_j + \psi) + \xi_j \right) \right) dF_{\psi,p}(\psi, \xi).$$

To solve for the Chebyshev weights, we use the Galerkin method described in Judd (1992). We define the residual function:

$$R(p; \gamma) = \gamma^{r'} \Lambda(p) - \dots \delta \int \log \left(\exp \left(\gamma^{r'} \Lambda(p\rho + \psi) \right) + \sum_j \exp \left(\beta_j^r - \alpha^r (p' \rho_j + \psi) + \xi_j \right) \right) dF_{\psi, \xi}(\psi, \xi) \quad (24)$$

Next, we let X be the matrix of K Chebyshev polynomials at each of the G points on our grid (i.e. G nodes). Our goal is to search for parameters, γ , that set the following expression to zero:

$$X' R(p; \gamma) = 0. \quad (25)$$

We use an iterated least squares approach for NFP.

1. Pick a starting value $\gamma^{r,0}$, $v_0^{r,0}(p; \Theta^r) = \gamma^{r,0'} \rho(p)$
2. Compute $Y(p; \gamma^{r,0}) = \delta \int \log \left(\exp \left(\gamma^{r,0'} \Lambda(p\rho + \psi) \right) + \sum_j \exp \left(\beta_j^r - \alpha^r (p' \rho_j + \psi) + \xi_j \right) \right) dF_{\psi, \xi}(\psi, \xi)$ using quadrature
3. solve the least squares problem: $\min_{\gamma} R(p; \gamma)' R(p; \gamma) \Rightarrow \min_{\gamma} (X\gamma^r - Y(p; \gamma^{r,0}))' ((X\gamma^r - Y(p; \gamma^{r,0})))$
 - for which the solution is: $\gamma^{r,1} = (X'X)^{-1} X'Y(p; \gamma^{r,0})$.
4. Compute $v_0^{r,1}(p; \Theta^r) = \gamma^{r,1'} \Lambda(p)$
5. Repeat steps 2 and 3 until convergence:.

D Jacobian of the Density of (p_t, S_t) in the Dynamic BLP model

The Jacobian is defined as follows:

$$J_{t,u \rightarrow Y} = \begin{bmatrix} \frac{\partial \psi_t}{\partial p_t} & 0 \\ 0 & \frac{\partial \xi_t}{\partial S_t} \end{bmatrix}.$$

Since $\frac{\partial \psi_t}{\partial \log(p_t)} = I_J$, we only need to compute the matrix of derivatives, $\left[\frac{\partial \xi_t}{\partial S_t} \right]$. We can simplify this calculation by applying the implicit function theorem to the following system

$$G(S_t, \xi_t) = s(p, \xi_t; \Theta) - S_t = 0$$

and computing the lower block of the Jacobian as

$$\begin{aligned} J_{t,\xi \rightarrow S} &= - \left[\frac{\partial G}{\partial \xi_t} \right]^{-1} \left[\frac{\partial G}{\partial S_t} \right], \\ &= \left[\frac{\partial s}{\partial \xi_t} \right]^{-1}, \end{aligned}$$

where the (j, k) element of $\frac{\partial s_{j,t}}{\partial \xi_{k,t}}$ is

$$\frac{\partial s_{j,t}}{\partial \xi_{k,t}} = \begin{cases} \sum_r \lambda_{r,t} s_j(p_t, \xi_t; \Theta^r) (1 - s_j(p_t, \xi_t; \Theta^r)) & , \text{ if } j = k \\ - \sum_r \lambda_{r,t} s_j(p_t, \xi_t; \Theta^r) s_k(p_t, \xi_t; \Theta^r) & , \text{ otherwise.} \end{cases}$$