

ECON8853 - Industrial Organization I

Lecture Notes from Julie H. Mortimer's lectures

Paul Anthony Sarkis
Boston College

Contents

1	Overview of Demand Systems and Vertical Model	6
1.1	Introduction	6
1.1.1	Intuition	6
1.1.2	An early example: Bresnahan (1987)	7
1.1.3	Approaches to demand estimation	7
1.2	Single product demand	8
1.2.1	Representative agent	8
1.2.2	Heterogeneous agents	8
1.3	Multi-product demand	9
1.3.1	Which approach to chose?	9
1.3.2	Product Space approach	10
1.3.3	Characteristic Space Approach	15
1.4	Application: Vertical model	18
1.4.1	Recovering the shares	19

1.4.2	Estimation	20
1.4.3	Issues with the vertical model	21
2	Logit, Nested Logit and Multinomial Probit	22
2.1	Identification of Choice Models	22
2.1.1	Only Differences in Utility Matter	22
2.1.2	Irrelevance of utility scale	23
2.2	“Plain-vanilla” Logit	24
2.2.1	Choice probabilities	24
2.2.2	Logit and taste variation	27
2.2.3	Derivatives and Elasticities	28
2.2.4	Independence of Irrelevant Alternatives (IIA)	29
2.2.5	Consumer Surplus	30
2.3	Nested Logit	31
2.3.1	IIA and substitution patterns	33
2.3.2	Research considerations	33
2.4	Multinomial Probit	34
2.4.1	Intuition	34
2.4.2	Choice probabilities	34
2.4.3	Applications	35
2.5	Estimation Strategy for Product-level data	35

2.5.1	Inversion in the logit case	35
2.5.2	Inversion in the nested-logit case	36
2.5.3	Endogeneity issues	37
2.5.4	Adding supply-side restrictions	38
3	Mixed Logit (BLP)	39
3.1	Model	40
3.1.1	Building the shares	41
3.2	BLP algorithm	42
3.2.1	Identification	45
3.2.2	Helping estimation	45
3.3	DFS algorithm	45
4	Product Availability	46
4.1	Introduction	46
4.2	Assortment variation	47
4.2.1	Unobserved assortment variation	47
4.2.2	Observed assortment variation	50
4.3	Stockouts	51
4.3.1	Partially observed stockouts	52
4.3.2	Observed stockouts	53
4.4	Limited consumer information	53

5	Entry Models	54
5.1	Introduction	54
5.1.1	Industry structure	54
5.1.2	Concentration measures	54
5.2	Sutton (1991)	54
5.2.1	Perfect competition	55
5.2.2	Imperfect competition	55
5.2.3	Endogenous sunk costs	56
5.3	Bresnahan and Reiss (1991)	56
5.4	Berry (1992)	59
5.5	Two-period models	61
5.6	Entry models with multiple equilibria	63
5.6.1	Entry games with structural errors	64
5.6.2	Entry games with expectational errors	67
6	Moment Inequalities	69
6.1	Framework	69
6.1.1	The agent's decision problem	69
6.1.2	Observables and disturbances	71
6.1.3	Moment inequalities	71
6.2	Applications	71

7	Single-agent Discrete Dynamic Programming	72
7.1	Definition	72
7.2	Rust (1987)	77
7.2.1	Setup	77
7.2.2	Econometric model	79
7.2.3	Estimation Method	82
7.3	Hotz-Miller approach	83
8	Retailing and Inventories	84
8.1	Aguirregabiria (1999)	84
8.1.1	Summary	84
8.1.2	Discussion	86
8.2	Hendel and Nevo (2002)	86
8.2.1	Summary	86

Chapter 1

Overview of Demand Systems and Vertical Model

1.1 Introduction

Demand system estimation is at the core of the Industrial Organization field of economics. In fact, it is central to many economic applications such as the study of comparative statics, welfare impacts, advertising, etc. Since there are many types of markets where estimating demand is useful, there are also many approaches to estimate demand. The goal of this chapter is to present the standard and most common approaches, work out the intuition to apply these methods and discuss their advantages and disadvantages.

1.1.1 Intuition

Before diving in the theory, let's get some intuition. A demand system is the relationship between prices and goods purchased on the consumer side. Thus, IO researchers typically want to estimate the effect of products characteristics, most commonly price, on consumers' propensity to buy the products. To do this, the IO researcher needs to observe a market and its interactions. A unit of observation is therefore a market interaction such as the purchase of a good. In

this unit of observation, we would ideally observe what good was purchased, its defining characteristics, its price, the state of available competing products, etc. This makes the data requirements pretty big, which is why acquiring a good enough dataset is complicated and often expensive.

1.1.2 An early example: Bresnahan (1987)

Demand system estimation is not only used for the sake of demand analysis, it can be a very useful tool to study broader industry topics. The use of demand models to talk about an industry as a whole is the main contribution of the NEIO movement (New Empirical IO). Bresnahan (1987) is one of the first empirical assessment of competition using demand estimation. The paper studies the presence of collusion in the automobile industry by looking at a dramatic price decrease event that happened in 1955.

The intuition of this paper is quite simple: assuming marginal costs stayed relatively constant during that period, Bresnahan estimates the variation in demand elasticities and compares it to a change in competition model. Using data on 85 models over 3 years, he finds that a change from a collusive to a more competitive equilibrium is consistent with the estimated change in elasticities.

1.1.3 Approaches to demand estimation

We consider different approaches to demand estimation based on three main divides:

1. Single product vs. multi-products: whether the system to estimate allows for differentiated products or a unique product.
2. Representative agent vs. heterogeneous agents: whether the system allows for consumers to have different tastes and behavior or not.
3. Product space vs. characteristic space: whether the system considers products as whole entities or as combinations of characteristics.

1.2 Single product demand

1.2.1 Representative agent

A demand system is generally thought to be a (more or less) general approximation of a true underlying demand curve, for example:

$$\ln(q) = \alpha p + X\beta + \varepsilon$$

where the ε captures any variation that is unobserved by the econometrician. Using that equation, the objective of the econometrician is to recover the underlying parameters (α, β) in order to understand the effects of prices and characteristics on demand. Note that this specification assumes a representative consumer choosing a (total) quantity q for a given price p . On the other side of the market, the firms that provide the good might also choose their prices p and characteristics x as a function of the expected demand, which is correlated with q . If this is the case, we say that we have endogeneity (from simultaneous equations, as we've seen in econometrics). Since the coverage of IV estimation in that case is more of an econometrics topic, we will assume knowledge of this in the rest of the chapter.

1.2.2 Heterogeneous agents

While the previous model is very simple, a straightforward extension would be to allow for heterogeneous agents on the market. To do that, we need to extend the previous model in two ways: (1) use a micro-founded model for individual demand that aggregates nicely and (2) estimate aggregated demand as:

$$\ln(q) = \int \gamma_i g(\gamma) d\gamma + \int \alpha_i p f(\alpha) d\alpha + \int \beta_i x h(\beta) d\beta + \varepsilon$$

where α, β and γ follow known distributions with unknown parameters to be estimated. This extension gives us some flexibility in terms of the potential interpretations. In a simple case where there are two types of agents: women and men for example, we will be able to estimate different tastes for different characteristics.

1.3 Multi-product demand

When we enter demand systems with multiple products, we need a way to differentiate those products. There are two approaches to do this: the product space approach, where a product is a nonseparable single entity; and the characteristic space approach, where a product is a combination of various characteristics.

1.3.1 Which approach to chose?

Disadvantages of the product space approach

In the product space approach, all competing products are estimated as fixed effects, meaning that, for a demand system with J products, each demand equation will have J parameters to estimate (at least). As the number of products increase, the number of parameters to estimate will become huge (J^2 in the simplest model). This issue is called the “too many parameters” problem. Note that some alternatives with groups of products will reduce this issue, while not solving it completely. A second issue with this approach is the introduction of a new good in the market, since it is not possible to estimate fixed effects when products are not yet available.

Disadvantages of the characteristic space approach

In the characteristic space approach, the data requirements are way higher than in the product space approach. This creates two main issues: first, obviously it will be hard to get all the data about all characteristics for all products; second, inevitably some characteristics will not be observed and might create endogeneity if they have an important role in demand. Moreover, we will get the same issue as in the product space approach when new dimensions are introduced (not new products, since we would be able to estimate consumers’ preferences using other products).

1.3.2 Product Space approach

This course explores product space approaches only briefly since most of modern IO revolves around the characteristic space approach.

As we've seen, the product space approach considers each product as an entity on its own, meaning that all observed and unobserved features of the product enter the "utility" derived from the product. Because different consumers can derive different utilities from the same product, we might be interested in using heterogeneous agents models. In fact, in this approach, the basic breakdown between models is whether they include a representative agent or heterogeneous agents. The heterogeneous agents models use more of the available information (better estimates) but also provide a framework for distributional policy analysis (estimating the distribution of responses). While these models have been around for a long time, they still face big issues such as treatment of zeroes (when some consumers have no access to a particular good) or aggregation issues (computational). This is why much of the work that has been done in product space revolves around representative agent models, starting with the AIDS model by Deaton and Muellbauer (1980).

In a representative agent model, the general demand model is defined as:

$$\ln q_j = \alpha p_j + \beta p_{-j} + \gamma x_j + \epsilon_j$$

where p_j represents the own-price, p_{-j} is the vector of all competing products' prices and x_j is a vector of all other observed characteristics.

As we know, prices in this model are endogenous (simultaneity issue) but here we have more than one endogenous price (we have as many as products), thus estimating this model will require a very demanding IV strategy. Moreover, we can see that each product's demand function has at least J parameters, so that the number of parameters to be estimated will increase quadratically compared to the number of products: we will typically need to reduce the dimensionality of the problem to be able to estimate the demand model.

Reducing the dimensionality

There are three popular approaches to reduce the dimensionality of the product space problem. The first two approaches are useful only for a particular set of questions, such as situations when variety is very important (e.g. in international trade they use them all the time), while the last one is the most often used in IO.

The first approach is to use a Constant Elasticity of Substitution model, such that utility is defined as:

$$u(q) = \left(\sum_k q_k^\rho \right)^{1-\rho} \Rightarrow q_j(p) = \frac{p_j^{-1/(1-\rho)}}{\sum_k p_k^{-\rho/(1-\rho)}} \cdot M$$

This approach is very efficient at reducing dimensionality since we now have only one parameter to estimate: ρ . However, it comes with the property that all goods have the same own-price elasticity and the same cross-price elasticities. While this definitely seems implausible, this approach is used in macro models and in trade.

The second approach is to use a logit demand model (not the same logit as in the characteristic approach). In this model, we have that:

$$u(q) = \sum_k (\delta_k - \ln(q_k)) q_k$$

The problem with this approach is the IIA property, which implies that elasticities depend on market shares rather than “closeness” between products.

Finally, the third empirical approach is to use a model of multi-level budgeting, as in a “utility tree”. This approach’s intuition is to divide the set of products into small groups and sub-groups, while allowing flexible substitution within groups. To use this approach, we need to assume two properties:

- Separability: this property ensures that preferences for products coming from one group are independent of consumption of products from other groups. Formally, this implies that the utility function takes the following form:

$$u(q) = f(u_1(q^1), \dots, u_n(q^n))$$

where there are n groups.

- Multi-stage budgeting: this property ensures that consumers can allocate expenditures in stages, considering only the current level of grouping and not what is inside the group.

It follows that this approach has three steps:

1. Grouping the products together in a (defensible) way.
2. Allocate expenditures to each group.
3. Allocate expenditures dynamics between groups.

Almost Ideal Demand System (AIDS)

One of the most popular model of the multi-level budgeting approach is the Almost Ideal Demand System (AIDS) model, designed by Deaton and Muellbauer.

Within a group g , demand for a product i is defined as:

$$w_i = \alpha_i + \beta_i \ln(y_g/P_g) + \sum_{j \in g} \gamma_{ij} \ln p_j + \epsilon_i$$

where w_i is total expenditure on product i , y_g/P_g is the real expenditure on group g and p_j is the price of other products in group g . As such, we can see that demand is not affected by products in other groups, except through the effect on the group g 's expenditure. Note that y_g/P_g depends on the price index of group g , which can be computed in multiple ways. The two most important ways to compute it are the simple logarithmic index:

$$P_g = \sum_{j \in g} w_j \ln p_j$$

which is a weighted average of log prices of products inside group g (weights are given by the expenditure level of each product). The second way is the exact price index of Deaton and Muellbauer (based on Shepard's lemma):

$$P_g = \alpha_0 + \sum_{j \in g} \alpha_j p_j + (1/2) \sum_{j \in g} \sum_{k \neq j \in g} \gamma_{jk} \ln p_j \ln p_k$$

which is more complex to estimate (more parameters).

Then, if there is a middle level between groups, demand can be estimated with AIDS again (using price indices instead of product prices) or by a simpler log-log specification. At this level of decision-making, both approaches are the same (in the sense that they do not follow a theoretical model). A log-log specification would look like:

$$\ln q_g = \alpha_g + \beta_g \ln y + \sum_h \delta_{g,h} \ln P_h + \epsilon_g$$

Finally, the top-level is a single log-log equation, with the addition of a set of demand shifters:

$$\ln q = \alpha + \beta \ln y + \delta \ln P + Z\gamma + \epsilon$$

Issues with AIDS

All in all, AIDS works very well when products are actually grouped in certain categories. This is the reason why it is widely used in the trade and macro-consumption literatures. It can be used in IO settings where groups could be single products but it requires instrumenting for every good!

Think of the case where there are J goods, then you need to estimate J^2 elasticities or at least $J \times (J + 1)/2$ if you are assuming that cross-price elasticities are symmetric. Usually, IO economists try to escape the problem by actually grouping the goods, however all the results will be dependent on your initial choice of grouping. Researchers have to think very deeply in the mechanism they use for grouping. Grouping subject to particular characteristics lead you closer to characteristic space models, but even then, how do you group when the characteristics are continuous etc.

Finally, let's look at one of the main challenges with the product space approach. Suppose a new good is introduced in a market, how can you use the previous demand estimated which did not account for the new product? This creates two issues:

1. Equilibrium effects: the introduction of a new good typically has significant effects on the pricing of the previously existing goods.
2. Extrapolation: using the previously estimated demand implies projecting demand where it is not defined.

Examples: Hausman, Leonard and Zona (1994)

The goal of the paper is to estimate demand for beer in the US in order to perform a merger analysis and test assumptions about the firms' conduct.

The market is divided hierarchically in three levels:

1. The top-level is beer against other goods. This is estimated using a log-log expenditure function:

$$\ln e_t = \beta_0 + \beta_1 \ln y_t + \beta_2 \ln P_{bt} + Z_t \delta + \varepsilon_3$$

where e_t is total expenditures, y_t is income and P_{bt} is the beer price index.

2. Within the beer group, the middle-level are the segments: premium, popular and light. Again, a log-log functional form is assumed so that for a given segment m , the demand for the segment is:

$$\ln q_m = \beta_m \ln e_t + \sum_{m'} \sigma_{m'} \ln \pi_{m'} + \alpha_{m'} + \varepsilon_2$$

3. Within each segment, the bottom-level contains five brands. This level is assumed to follow an AIDS:

$$w_k = a_k + \sum_j a_{jk} \ln p_j + \beta_k \ln(x/P) + \varepsilon_1$$

In the bottom-level, you have to instrument the price since it is correlated with unobserved product quality (taste, etc.) and unobserved demand shocks (special events like the World Cup, etc.). Then you must find brand-level instruments: this is where the infamous Hausman instruments come in. Hausman instruments use prices in one city to instrument for prices in other cities. These instruments are typically strong, however, their relevance is very questionable. In particular, think about nation-wide advertising which was the main criticism of these instruments. They would clearly affect both cities' prices in the same way creating correlation. Most of the time you can observe these patterns but be wary of unobserved shocks that could affect the complete market. Optimally, the best instrument would be observed input costs for the brand, tax rates, etc.

1.3.3 Characteristic Space Approach

The characteristic space approach follows a different philosophy than the product space:

- Products are considered as a vector of characteristics
- Consumers' preferences (and thus their utility functions) are defined over these characteristics.
- Consumers are assumed to choose the characteristic vector (the product) that maximize their utility.
 - One consumer makes one choice to buy a single product. Allowing for buying more than one product is computationally costly and is an open area for research.
- Demand is aggregated by simply summing over consumers' choices.

Formal base model

Assume the following context:

- A consumer i is offered J alternatives. Consumer i is identified by his characteristics, v_i , similarly to the product being defined by its characteristics x_j and its price p_j .
- He must choose one option only $j \in J$, which he will do with probability P_{ij}
- Utility of an individual i for a good j is given by:

$$U_{i,j} \equiv U(x_j, p_j, v_i, \theta) \text{ for all } j = 0, 1, \dots, J$$

where good $j = 0$ is referred to as the outside good.

The variable (or most probably the vector) x_j contains non-price characteristics about good j , for example the size of the engine, the AC system, the dashboard software, etc.; while v_i is the vector of consumer characteristics such as income, age, etc. Finally, θ is the vector of parameters of the utility which is to be estimated. Note that in order to estimate the model correctly, you need to be sure that all consumers (or groups of consumers) in the sample have the same access to all J goods, or in other words, they face the same choice set.

Consumer i will choose good j if and only if $U_{i,j} > U_{i,k}$ for any other good k . This means that the probability of consumer i buying good j is given by:

$$\Pr [U_{ij} > U_{ik} \text{ for all } k \neq j]$$

Now assuming that utility is given by an observable component V_{ij} and an unobservable component ε_{ij} , we can rewrite the probability of buying the good as:

$$\begin{aligned} P_{ij} &= \Pr [U_{ij} > U_{ik} \text{ for all } k \neq j] \\ &= \Pr [V_{ij} + \varepsilon_{ij} > V_{ik} + \varepsilon_{ik} \text{ for all } k \neq j] \\ &= \Pr [\varepsilon_{ij} - \varepsilon_{ik} > V_{ik} - V_{ij} \text{ for all } k \neq j] \end{aligned}$$

Finally, define $f(\varepsilon_i)$ as the J -dimensional probability density function for consumer i . Then the probability P_{ij} can be written as:

$$P_{ij} = \int \mathbb{I}\{\varepsilon_{ij} - \varepsilon_{ik} > V_{ik} - V_{ij}\} f(\varepsilon_i) d\varepsilon_i$$

which is (as ε_i) a J dimensional integral over $f(\varepsilon_i)$.

We will see later on that some choices could make our life easier in estimating this integral (hard to do it analytically), mainly by parameterizing the error term:

- **Multivariate normal:** $\varepsilon_i \sim N(0, \Omega)$ \rightarrow multinomial probit model.
- **Type 1 Extreme-value:** $f(\varepsilon_i) = e^{-\varepsilon_{ij}} e^{-e^{-\varepsilon_{ij}}}$ \rightarrow multinomial logit model.
- other variants...

We could also turn to thinking about absolute demand for good j instead of its probability, which is given by the size of the set $S_j(\theta)$ defined as:

$$S_j(\theta) \equiv \{i : U_{i,j} > U_{i,k} \text{ for all } k\}$$

Now suppose consumers are distributed following a pdf $f(v|\theta)$, we could recover the market share of good j as:

$$s_j(\mathbf{x}, \mathbf{p}|\theta) = \int_{i \in S_j(\theta)} f(v) dv$$

Given a total market size of M , total demand for good j is: $M \cdot s_j(\mathbf{x}, \mathbf{p}|\theta)$

As we've seen the complete model of characteristic space relies on utility functions and error terms distributions. This is why, unsurprisingly, the functional forms of utility and errors are exactly what will differentiate the models within the characteristic space approach.

There are 5 main models of characteristic space.

Pure horizontal model

Analog to the Hotelling model (in its simplest form), there are n ice cream sellers along a beach where consumers are distributed. A consumer i has utility:

$$U_{ij} = \bar{u} - p_j + \theta(\delta_j - v_i)^2$$

where the term $(\delta_j - v_i)^2$ captures some kind of quadratic transportation cost in the distance from i 's preferences (v_i) to the characteristics of product j (δ_j). In the simple Hotelling model, δ_j and v_i are location parameters (where ice cream seller j and consumer i are on the beach).

Pure vertical model

On the pure vertical, you don't need to model a distance from consumer preferences since everyone agrees on what product is better. In this case, utility is given by:

$$U_{ij} = \bar{u} - v_i p_j + \delta_j$$

The interaction between consumer characteristics and prices show that although every consumer has the same utility for the characteristics of good j (i.e. $\bar{u} + \delta_j$), consumers differ in their willingness to pay (v_i).

Logit model

In the logit model, consumers have the same taste for the goods' characteristics but are subject to an idiosyncratic shock depending on both the product and the

consumer. Utility is given by:

$$U_{ij} = \delta_j - p_j + \varepsilon_{ij}$$

where ε_{ij} follows an iid extreme-value type I distribution. Analogous to the OLS, the error term allows for identification of the other parameters as long as it is not correlated with the δ_j and the price.

Pure characteristic model

This model nests both concepts of verticality and horizontality of differentiation.

Utility is given by:

$$U_{ij} = f(v_i, p_j) + \sum_k \sum_r g(x_{jk}, v_{ir}, \theta_{kr})$$

so that it could handle vertical model if $g(\cdot)$ is 0 for cross-products variables. It could also handle the horizontal model. It is a very general form and rarely used?

BLP model

The BLP model (for the authors' names Berry, Levinsohn and Pakes), is a parameterized version of the pure characteristics model. It is probably the most commonly used demand model in the empirical literature as of now, and is also very popular in non-research applications of IO. Utility is given by:

$$U_{ij} = f(v_i, p_j) + \sum_k \sum_r x_{jk} v_{ir} \theta_{kr} + \varepsilon_{ij}$$

where the ε_{ij} extreme-value Type I error term has been added.

1.4 Application: Vertical model

As we have mentioned in earlier, in a vertical model, all consumers agree on the relative quality of products (i.e. there is a clear ranking). However, people differ

in their willingness to pay for such quality. The general form of utility is:

$$U_{ij} = \bar{u} - v_i p_j + \delta_j$$

From this equation, the objective is to recover the implied shares as a function of the parameters, and then estimate these parameters by comparing with the observed shares in the data.

1.4.1 Recovering the shares

In order to recover the shares, we first need to know the shape of the distribution of v_i . In this application, we assume it follows a lognormal distribution with implicit mean μ and variance σ^2 .

Then, we order all goods in increasing order of price (or δ_j) but by assumption, it has to be the same order since people would buy a better alternative if it was also cheaper.

We can show that if $U_{ij} < U_{ij+1}$, then, $U_{ij} < U_{ik}$ for all $k \geq j + 1$, assuming that $(\delta_{j+1} - \delta_j)/(p_{j+1} - p_j)$ is decreasing in j (the improvement in quality is not as fast as the increase in prices).

Normalize the outside good to 0, and you get that a consumer would choose the outside good if and only if:

$$0 \geq \max_{j \geq 1} U_{ij}$$

Using the property of the previous paragraph, we know that this is equivalent to:

$$0 \geq U_{i1}$$

Thus, we have that the probability that a consumer buys the outside good is:

$$P_{i0} = \Pr[0 \geq -v_i p_1 + \delta_1] = \Pr[v_i \geq \delta_1/p_1]$$

and more generally, the share of the outside good is given by:

$$s_0(\theta) = 1 - F(\delta_1/p_1)$$

where θ is the parameter vector containing all δ s, μ and θ and $F(\cdot)$ is the cdf of the lognormal distribution.

For inside goods, we have that a consumer would choose good 1 if and only if $U_{i1} > 0$ and $U_{i1} > U_{i2}$. Thus, the probability is:

$$\begin{aligned} P_{i1} &= \Pr [0 \leq -v_i p_1 + \delta_1 \geq -v_i p_2 + \delta_2] \\ &= \Pr [v_i \leq \delta_1/p_1 \text{ and } v_i \geq (\delta_2 - \delta_1)/(p_2 - p_1)] \end{aligned}$$

and more generally, the share of the first good is:

$$s_1(\theta) = F(\delta_1/p_1) - F((\delta_2 - \delta_1)/(p_2 - p_1))$$

Following the same argument for the following products, we get:

$$\begin{aligned} s_j(\theta) &= F[(\delta_j - \delta_{j-1})/(p_j - p_{j-1})] - F[(\delta_{j+1} - \delta_j)/(p_{j+1} - p_j)] \text{ for all } j = 2, \dots, J-1 \\ &\text{and } s_J(\theta) = F[(\delta_J - \delta_{J-1})/(p_J - p_{J-1})] \end{aligned}$$

This gives us a system of J equations with $J + 2$ unknown parameters in θ . This means that we can only identify J parameters (the δ s, provided we know/assume the parameters of the consumer characteristic distribution). Each δ will be identified using the associated market share and price.

1.4.2 Estimation

As we have seen in the previous section, our pure vertical model is not identified if the parameters of v_i 's distribution are unknown. That's why, in this particular case, we need to reduce the dimension of the problem first, and then estimate via maximum likelihood.

In the vertical model, we need to “project” the product quality δ_j onto a few characteristics x_j (a K dimension vector, such that $K < J$). Formally, we have that:

$$\delta_j = \sum_{k=1}^K \beta_k x_{kj}$$

Now that all goods have the same K coefficients governing the value of δ , θ is effectively a $K + 2$ dimension vector, estimable using the J shares!

1.4.3 Issues with the vertical model

There is a number of issues that arise when using this model:

1. Cross-price elasticities are 0 for all goods that are further away than direct neighboring products. This cannot capture a lot of actual variation in these elasticities.
2. Own-price elasticities are not decreasing in j , thus leading to high-priced products having comparable elasticities to low-priced products, which is thought of as an undesirable property.
3. Estimating the model relies on multiple researcher assumptions:
 - The outside good has to be chosen to be the lowest or highest good.
 - The functional form of the distribution of v_i has to be assumed.
4. The error term is deterministic given the functional form of $f(v)$.

Chapter 2

Logit, Nested Logit and Multinomial Probit

2.1 Identification of Choice Models

In order to identify all choice models we are going to see in this chapter, we will need to make several specification assumptions. These assumptions will allow us to go around two main issues that arise because of the choice process. These issues can be summarized in two statements: “only differences in utility matter” and “the scale of utility is arbitrary”.

2.1.1 Only Differences in Utility Matter

Recall the choice process between products as described in the previous chapter: a consumer i will choose a product j over its alternatives if the utility derived from j , denoted U_{ij} is greater than of all other alternatives. This means that the consumer will choose the product with the greatest perceived utility, and adding the same constant to all alternatives will not change the decision. This phenomenon has a particularly important implication: only parameters that capture differences across alternatives can be identified.

Product-specific constants

It is often reasonable to include an unobserved product-specific constant in the utility derived from a product j such that $U_{ij} = V_{ij} + \varepsilon_{ij}$, where

$$V_{ij} = X_j\beta + \xi_j$$

However, since only differences in utility matter, we will only be able to capture difference between ξ s, not their absolute level. To see that consider a binary choice model where $\tilde{\xi}_j - \tilde{\xi}_k = \xi_j - \xi_k$, then, we would have

$$U_{ij} > U_{ik} \Leftrightarrow \tilde{U}_{ij} > \tilde{U}_{ik}$$

and we would not be able to identify which value of ξ or $\tilde{\xi}$ is the right one. That is why, in a model with J alternatives, at most $J - 1$ alternative specific constant can enter the model, and one product must have its ξ normalized to zero.

2.1.2 Irrelevance of utility scale

Just as adding a constant to the utility of all alternatives does not change the decision maker's choice, neither does multiplying each alternative's utility by a constant, i.e. the product with the highest utility would stay the same. The scale of utility and the variance of the error term are linked by definition, since multiplying U_{ij} by a constant λ for all j would imply multiplying the variance of ε_{ij} by λ^2 for all j . That is why the econometrician will have to normalize the model by normalizing the variance of the error term.

Normalization with iid errors

In the case of iid error terms, normalizing the variance is straightforward since normalizing the variance of one error term will normalize the variance of all error terms.

Normalization with heteroskedastic errors

In the case of heteroskedastic error terms, when at least one segment of the population/alternatives has a different variance than other segments, normalizing the model has to be done by normalizing the variance of one segment and estimating the variance for all others relative to that segment.

Normalization with correlated errors

Finally, when the error term is correlated across products, normalizing for scale is more of an issue, indeed, in that case, normalizing the variance for one product is not enough to normalize the variance of other products. The solution is to set the variance of one of the error differences to a number, then estimating all other variances and covariances relative to the constant.

2.2 “Plain-vanilla” Logit

2.2.1 Choice probabilities

To derive the logit model, recall the notation used in the previous chapter where each product has a mean quality level δ_j that incorporates the characteristics of the product as well as its price, and an idiosyncratic error term ε_{ij} . Thus, we have:

$$U_{ij} = \delta_j + \varepsilon_{ij}$$

In the logit model, we assume that the error term is drawn iid, from a Type-I Extreme Value (or Gumbel) distribution. The Gumbel distribution has the following laws of probability:

$$f(x) = e^{-x} e^{-e^{-x}} \text{ and } F(x) = e^{-e^{-x}}$$

Moreover, the variance of the distribution is assumed to be $\pi^2/6$; this assumption is crucial for identification since it implies normalization of the scale of utility (this is explained in the first section of this chapter). The mean of this distribution

is not 0 (!) but since only differences between utility levels are important, this fact will not cause any issues. Finally, this distribution assumption allows us to have a closed-form solution for the distribution of the difference between two error terms: specifically, let $\tilde{\varepsilon}_{ijk} = \varepsilon_{ij} - \varepsilon_{ik}$, we have that:

$$F(\tilde{\varepsilon}_{ijk}) = \frac{\exp(\tilde{\varepsilon}_{ijk})}{1 + \exp(\tilde{\varepsilon}_{ijk})}$$

This distribution law is almost exactly the same as the Normal distribution, so that even if the extreme-value assumption gives fatter tails and slight asymmetry, we get back a Normal distribution in the differences. Nevertheless, note that the shape of the distribution is not that important and that the more important assumption is independence of errors, which requires a well-specified model to be justifiable (i.e. no missing variables, etc.). The independence assumption also requires that the unobserved part of utility for each product (ε_{ij}) is completely independent with respect to other products! In other words, taste for products cannot be correlated in an unobserved way. If this is the case, then we need to use other models that will be described later.

The probability that consumer i buys good j , given ε_{ij} , is the conditional probability that he chooses j over all other k :

$$\begin{aligned} P_{ij}|\varepsilon_{ij} &= \Pr [U_{ij} > U_{ik} \text{ for all } j \neq k | \varepsilon_{ij}] = \Pr [\delta_j + \varepsilon_{ij} > \delta_k + \varepsilon_{ik} \text{ for all } j \neq k] \\ &= \Pr [\varepsilon_{ik} < \varepsilon_{ij} + \delta_j - \delta_k \text{ for all } j \neq k] \\ &= \prod_{j \neq k} F(\varepsilon_{ij} + \delta_j - \delta_k) \text{ (by indep.)} \end{aligned}$$

We can then proceed to compute the unconditional probability by integrating over all possible values of ε_{ij} from the Gumbel distribution:

$$P_{ij} = \int \left(\prod_{k \neq j} F(\varepsilon_{ij} + \delta_j - \delta_k) \right) f(\varepsilon_{ij}) d\varepsilon_{ij} = \frac{e^{\delta_j}}{\sum_{k=0}^J e^{\delta_k}}$$

which is called the logit choice probability, and also equivalent to the market share of product j , under the assumption that all consumers are identical in their distributions.

Some properties

This model is very practical in the sense it is very easy to set up and displays several desirable properties:

- $s_j = P_{ij} \in (0, 1)$ so that any market share can be rationalized by the model, except in extreme cases with only one product having all the sales (in those cases we would not want a logit model anyway), or some products having no sales (then we should take them out of the sample).
- As δ_j increases, reflecting higher attributed quality, the market share will increase to 1. If δ_j decreases, then the market share goes to 0. However, note that these results only hold at the limit, meaning a good will never have a market share equal to 0 or 1. If that is the case in the data, the researcher should take the product out of the dataset. This is why estimating a logit model must be done on a time-period long enough so that all goods considered have been sold at least once (think about vending machines).
- $\sum_j s_j = 1$, meaning that one alternative will be chosen, always.
- $s_j = P_{ij}$ is everywhere continuously differentiable in the characteristics δ_j and in price p_j .

The last point is useful to compute price derivatives (and elasticities):

$$\begin{aligned} \frac{\partial s_j(\theta)}{\partial p_j} &= \frac{-\alpha e_j^\delta \cdot (1 + \sum_k e_k^\delta) + \alpha e_j^\delta e_j^\delta}{(1 + \sum_k e_k^\delta)^2} = \frac{-\alpha e_j^\delta \cdot [1 + \sum_k e_k^\delta - e_j^\delta]}{(1 + \sum_k e_k^\delta)^2} \\ &= -\alpha \cdot \frac{e_j^\delta}{1 + \sum_k e_k^\delta} \cdot \frac{1 + \sum_k e_k^\delta - e_j^\delta}{1 + \sum_k e_k^\delta} \\ &= -\alpha s_j(1 - s_j) \end{aligned}$$

And using the same method, we find $\partial s_j(\theta)/\partial p_k = \alpha s_j s_k$

Characterizing mean utility

The estimation of the mean quality characteristic δ_j can be done over a set of non-price characteristics X_j . In particular, consider the following:

$$\delta_j = X_j \beta - \alpha p_j + \xi_j$$

where α and β are taste parameters for the given characteristics and ξ_j is an independent, unobserved mean utility shock. This shock is important in order to rationalize any pattern of market shares (suppose a product has better characteristics in every way but still has lower shares).

As introduced in the first section of this chapter, product-specific constants can only be estimated by normalizing the ξ of one product to a constant. Here, the choice is already made since the outside good has a zero utility by assumption. However, another issue arises by introducing this term. If consumers know ξ_j , firms must also know it and price their products accordingly. This introduces endogeneity in the error term with respect to prices. Thus, the econometrician needs instruments to correct for that endogeneity. Typically, these instruments include cost shifters that would not affect intrinsic demand for the good but would definitely have an effect on prices. Note that any non-price characteristic included in X_j that would also be correlated with the unobservable characteristics ξ_j needs to be instrumented as well.

2.2.2 Logit and taste variation

The logit model is very efficient and easy to set up if you need to measure systematic taste variation, that is, variation in tastes associated with observable characteristics X_j . Variations associated with idiosyncratic randomness, ε_{ij} are assumed away as type-I extreme value random variables.

In reality, the value that agents put on a particular attribute varies across these individuals. Sometimes these tastes vary in an identifiable way, for example, low-income agents might be more concerned about the price than others; sometimes you might observe people with the exact same observable characteristics choosing different options. Logit models allow for identification of these taste variations only within limits: if this variation is linked with observable characteristics, then logit models can be applied; if this variation is purely random, then we will require other models.

As an example, consider the choice of households between cars, where two characteristics enter the decision: price p_j and shoulder room x_j . A household i

places value U_{ij} on buying good j such that:

$$U_{ij} = \beta_i x_j - \alpha_i p_j + \varepsilon_{ij}$$

where you can see that the parameters α, β vary across households i . Say the shoulder room taste β_i depends only on the number of members in the household m_i such that $\beta_i = \rho M_i$; this means that m_i is positively correlated with β_i . Similarly let's say the importance of price is negatively correlated to income I so that $\alpha_i = \theta/I_i$. Substituting back into the utility function we get:

$$U_{ij} = \rho M_i x_j - \theta p_j / I_i + \varepsilon_{ij}$$

which can be estimated with a standard logit model where variables are interaction between characteristics of the product and characteristics of the household.

In contrast, suppose taste variation was subject to an error term such that $\beta_i = \rho M_i + \mu_i$ where μ_i is non-observable (a random variable). Then the utility would be written as:

$$U_{ij} = \rho M_i x_j - \theta p_j / I_i + \varepsilon_{ij} + \mu_i x_j$$

and the logit would confound the error term with $\tilde{\varepsilon}_{ij} = \varepsilon_{ij} + \mu_i x_j$ which is clearly correlated across product specifications and households.

Bottomline is logit models can handle systematic taste variation but not random taste variation. For the latter, we will need more complex models, such as the BLP model studied in the next chapter.

2.2.3 Derivatives and Elasticities

The logit model presented above displays features that are not desired in the context of elasticities. These problems imply that one would run into issues while trying to expand logit results in terms of welfare analysis, antitrust analysis, etc.

First, recall the own- and cross-price derivatives and elasticities:

$$\frac{\partial s_j}{\partial p_j} = -\alpha s_j(1 - s_j); \quad \frac{\partial s_j}{\partial p_k} = \alpha s_j s_k; \quad \epsilon_{jj} = -\alpha p_j(1 - s_j); \quad \epsilon_{jk} = \alpha p_k s_k$$

These imply that:

- Two products with the same market shares will have the same markups. Indeed, from the first-order conditions, we have that:

$$p_j - mc_j = \frac{1}{|\epsilon_{jj}|} \Leftrightarrow \frac{p_j - mc_j}{p_j} = \frac{1}{\alpha(1 - s_j)}$$

which would be the same for $s_j = s_k$. This is obviously not intuitive but also not observed in reality.

- Own-price elasticities are higher for higher priced goods. This fact comes directly from the formula of the own-price elasticity that shows a positive effect (in absolute value) from the prices. This is counter-intuitive since it would imply that people buying higher-priced items would be more price-sensitive than people buying lower-priced goods.
- Substitution between goods only depends on relative shares and not proximity of product characteristics. In fact, cross-price elasticities are given by $\frac{\partial s_j}{\partial p_k} \frac{p_k}{s_j} = \alpha p_k s_k$, which is not a function of characteristics of neither products.

2.2.4 Independence of Irrelevant Alternatives (IIA)

We have seen in the overview of the logit model that idiosyncratic errors are independently distributed across products, following a type-I extreme value distribution. The choice of the distribution turns out not to be very important, but the independence property has some implications that the researcher should know about. In fact, the independence assumption means that the unobserved utility derived from one good (ε_{ij}) is unrelated unconditionally to the unobserved utility from another good. This assumption seems to be rather restrictive in the context of goods but helped us a lot in coming up with the solution of the logit model.

The independence of the error terms turns out to have problematic implications in the realism of the choice mechanism. In fact, suppose a consumer has a choice of going to work by car (c) or taking a blue bus (bb). Say that utility derived from both models are the same, such that $P_c = P_{bb} = 1/2$, meaning the ratio of probabilities is one. Now suppose that a red bus (rb) is introduced in the market such that it is identical in every aspect except the color to the blue bus. The

probability of taking the red bus should therefore be the same as taking the blue bus: $P_{rb}/P_{bb} = 1$. However, the ratio of probabilities between the car and the blue bus has not changed because of independence of irrelevant alternatives, meaning that $P_{bb}/P_c = 1$, thus $P_c = P_{bb} = P_{rb} = 1/3$ is the prediction of the logit model. In real life though, we would expect the probability to take the car to remain exactly the same since the actual problem is to either take the bus (regardless of the color) or the car, yielding $P_c = 1/2, P_{bb} = P_{rb} = 1/4$.

This IIA problem is nonetheless also a feature of the logit model when it corresponds to reality. In fact, this property allow the researcher to consider only subsets of the complete set of alternatives and still get consistent estimates, as long as for each observation, the actual choice is kept in the set. Another practical use of that property is that if the researcher is only interested in a few choices, then they do not need to include the other choices in the dataset, leaving the data research part out of the picture.

2.2.5 Consumer Surplus

For policy analysis, the researcher is often interested in measuring the change in consumer surplus that is associated with a particular event (introduction of a product, merger, etc.). Under the logit assumption, this value of the consumer surplus also takes a simple closed form that can be easily computed from the model. We know that each consumer will choose the product that yields the highest utility. In the aggregate case, all consumers have the same tastes so that in expectation (the econometrician does not observe ε_{ij}), consumer surplus is defined as the value of utility derived from the best good. Formally,

$$E[CS] = E \left[\max_j \{U_{ij}\} \right] = E \left[\max_j \{\delta_j + \varepsilon_{ij}\} \right]$$

which yields:

$$E[CS] = \ln \left(\sum_{j=1}^J e^{\delta_j} \right) + C$$

where C is an unknown that represents the fact that absolute utility cannot be measured. This value is called the logit inclusive value, or the log-sum term. As you can see, comparing two policies is easy in this setting, let one policy be

denoted by the superscript 0 and the other by the superscript 1. Then,

$$\Delta E[CS] = \ln \left(\sum_{\mathcal{J}^1} e^{\delta_j^1} \right) - \ln \left(\sum_{\mathcal{J}^0} e^{\delta_j^0} \right)$$

If we had transaction-level data with individual characteristics, we could perform this analysis of consumer surplus at each unit of observation and aggregate to find our results.

2.3 Nested Logit

A natural extension of the simple logit model, allowing for richer substitution patterns and a somewhat less restrictive IIA property is the nested logit model.

A nested logit model partitions the choice set in different subsets called nests such that the actual choice of one good follows from choices among nests. For example, in the simplest nested logit, with one nest, consumers first choose a nest (a category of products) as modelled in a simple logit, then within a nest, they choose a product inside the nest (again modelled as a simple logit). This sequence of logit models creates a different type of model called the nested logit. These nests are chosen by the researcher (which is not ideal) and they provide a strong structure to the model which could affect results significantly. Following the notation of Cardell (1991) and Berry (1994), we model utility as:

$$U_{ij} = \delta_j + \zeta_{ig} + (1 - \sigma)\varepsilon_{ij}$$

where the new terms are: ζ_{ig} an idiosyncratic “nest” taste shock that applies to all goods j in the nest g ; and σ a parameter of correlation in tastes within nests (if σ is high, tastes within group are very correlated and the nest structure matters, if σ is low, then the correlation in tastes within the nest is small and the nest structure is irrelevant). At first, it might seem that the σ term will not allow for the simple logit model, but in fact, the ζ_{ig} error term follows a unique distribution distribution function such that $\zeta_{ig} + (1 - \sigma)\varepsilon_{ij}$ follows an extreme-value distribution (not type-I however, we call it Generalized Extreme Value or GEV). This makes σ important to both terms, and in particular, when $\sigma \rightarrow 0$ we go to the simple logit, and $\sigma \rightarrow 1$ yields to more within-group correlation.

We know that within the same nest, the utility level $u_{ij} = \delta_j + \zeta_{ig}(1 - \sigma) + \varepsilon_{ij}$, can be reduced down to ignore the effect of ζ , since it is the same across products. We end up in the same setting as the simple logit model. This implies that the conditional share of product j , or the share within the nest, is given by the same formula as the simple logit:

$$s_{j|g} = \frac{\exp(\delta_j/(1 - \sigma))}{\sum_{k \in \mathcal{J}_g} \exp(\delta_k/(1 - \sigma))}$$

where the denominator for this expression can be written as D_g , the total demand for products in the group g .

Meanwhile, across groups, we have that both error terms still give a type-I EV, so again, we can think about the demand for a group as we did in the logit case:

$$s_g = \frac{D_g^{(1-\sigma)}}{\sum_h D_h^{(1-\sigma)}}$$

Finally, the share of a product j is given by the product of both the share of the group containing j and the conditional share of j within the group:

$$s_j = s_{j|g} \cdot s_g = \frac{\exp(\delta_j/(1 - \sigma))}{D_g^\sigma \cdot \left(\sum_h D_h^{(1-\sigma)} \right)}$$

As we can see, demand for product j depends on its own quality level relative to its group, the quality of the group relative to the other groups and σ , the correlation in tastes within nests.

The outside good in this model is considered as one of its own group, so that $s_{0|0} = 1$ and with the normalization of $\delta_0 = 0$ and $D_0 = 1$, we get:

$$s_0 = 1 \cdot s_0 = \frac{1}{\sum_h D_h^{(1-\sigma)}}$$

This analytical derivation is the same when we extend the model to have nests within nests and more.

2.3.1 IIA and substitution patterns

The nest structure of this model satisfies two properties:

1. For any two alternatives that are in the same nest, the ratio of probabilities is independent of the attributes and/or existence of other products in the nest or in other nests.

$$s_j/s_k = s_{j|g}/s_{k|g} = \frac{\exp(\delta_j/(1-\sigma))}{\exp(\delta_k/(1-\sigma))}$$

Equivalently, we say that IIA holds within the nest.

2. For any two alternatives that are in different nests, the ratio of probabilities depends on the attributes and/or existence of other products in the two nests.

$$s_j/s_k = s_{j|g}/s_{k|h} \cdot s_g/s_h = \frac{\exp(\delta_j/(1-\sigma))}{\exp(\delta_k/(1-\sigma))} \cdot \frac{D_g^{1-\sigma}}{D_h^{1-\sigma}}$$

Equivalently, we say that IIA does not hold across nests.

2.3.2 Research considerations

As we have seen, the nested logit is a fairly simple extension on the simple logit case, which relies on relaxation of the IIA property across groups (or nests). The model is often told in a narrative of sequential choice: consumers first choose a group g (s_g), then a product (or group) j within group g ($s_{j|g}$).

The results of the model depend very strongly on the ex ante nesting structure chosen by the researcher. Therefore, it is very important (and hard) to understand what the appropriate structure should be. The sequential narrative is supposed to help the researcher come up with an accurate structure but it might be unhelpful at times.

Identification comes in different flavors:

- Parameters associated with characteristics and prices are identified within group by variation in these exact characteristics.
- The correlation in tastes parameter (σ) is identified by variation in

2.4 Multinomial Probit

2.4.1 Intuition

The two previous models covered in this chapter had two main issues: they cannot handle random taste variation and they imply restrictive substitution patterns due to IIA (even though the nested logit solves this issue in some way). The multinomial probit on the other side, is not affected by these two issues at all. In essence, the MNP approach allows for correlation between choices thanks to a flexible covariance matrix of the errors.

2.4.2 Choice probabilities

As in the general case, we start with writing the individual (unconditional) probability of choosing product j :

$$\begin{aligned} P_{ij} &= E \left[\Pr \left[U_{ij} > U_{ik} \text{ for all } j \neq k | \varepsilon_{ij} \right] \right] \\ &= E \left[\Pr \left[\delta_j + \varepsilon_{ij} > \delta_k + \varepsilon_{ik} \text{ for all } j \neq k | \varepsilon_{ij} \right] \right] \\ &= E \left[\Pr \left[\varepsilon_{ik} < \varepsilon_{ij} + \delta_j - \delta_k \text{ for all } j \neq k | \varepsilon_{ij} \right] \right] \\ &= \int \prod_{j \neq k} F(\varepsilon_{ij} + \delta_j - \delta_k) f(\varepsilon_{ij}) d\varepsilon_{ij} \end{aligned}$$

But in the multinomial probit case, the law of probability for ε_{ij} depends on all the other error terms! Thus we have correlation between choices, that we estimate.

Formally, we have:

$$\varepsilon_i \sim N(\mu, \Omega)$$

where ε_i is the vector of all product error terms. This means that Ω is a $(J + 1)^2$ -elements symmetric matrix. In practice, we restrict μ to be 0 (especially if we estimate ξ s), but following the first section of this chapter we will also need to restrict the variance of the error terms heavily. In fact, as mentioned in that section, in the case of iid but correlated error terms, we need to normalize one diagonal term of the error-differences distribution and estimate the other terms relative to that term.

2.4.3 Applications

The MNP approach is interesting in most cases where choice correlations are important, but is limited by the fact that data requirements get huge when we increase the number of products.

2.5 Estimation Strategy for Product-level data

The general estimation strategy that applies to logit models with product-level data is detailed in the following steps:

1. Assume that data is drawn from markets with large n .
2. Assume that observed market shares are measured without errors.
3. For each θ (vector of parameters), there exists a unique ξ such that the model shares and observed shares are equal (J equations, J unknowns).
4. We invert the model to get ξ as a function of the parameters. How this step is performed depends on the functional form of the model.
5. Using ξ , we can create the moments of the model, estimating them by GMM.

2.5.1 Inversion in the logit case

We need to express the quality unobservable term ξ_j in terms of all observable characteristics. First, recall that:

$$\delta_j = X_j\beta - \alpha p_j + \xi_j$$

but in this equation, δ_j is unknown (there is no such thing as a perceived mean quality level). Nevertheless, we can use the market shares formulas that link the δ_j to the observed market shares s_j (observed without errors). Thus,

$$s_j = \frac{\exp(\delta_j)}{1 + \sum_k \exp(\delta_k)}$$

but again, we run into a problem since it is expressed as a function of other δ_k : we could solve a system of equations, or simplify the model using the normalized

good. Indeed,

$$\frac{s_j}{s_0} = \exp(\delta_j) \Leftrightarrow \ln(s_j) - \ln(s_0) = \delta_j \Leftrightarrow \ln(s_j) - \ln(s_0) = X_j\beta - \alpha p_j + \xi_j$$

which in turn yields the inversion equation:

$$\ln(s_j) - \ln(s_0) + \alpha p_j - X_j\beta = \xi_j$$

2.5.2 Inversion in the nested-logit case

In the same way as in the simple logit model, our main objective is to get the link between δ_j and observables so that we can identify ξ_j . In this case:

$$s_j = \frac{\exp(\delta_j/(1-\sigma))}{D_g^\sigma \cdot \left(\sum_h D_h^{(1-\sigma)}\right)} \text{ and } s_0 = \frac{1}{\sum_h D_h^{(1-\sigma)}}$$

so that

$$\ln(s_j) - \ln(s_0) = \frac{1}{1-\sigma} \cdot \delta_j - \sigma \ln(D_g)$$

We turn into a new problem caused by the two-level structure, which is D_g is not parameterized. Again, we solve this issue by using the normalized good:

$$s_g/s_0 = D_g^{(1-\sigma)} \Leftrightarrow \ln(s_g) - \ln(s_0) = (1-\sigma) \cdot \ln(D_g)$$

which we can plug back in the previous equation to get:

$$\begin{aligned} \ln(s_j) - \ln(s_0) &= \frac{1}{1-\sigma} \cdot \delta_j - \frac{\sigma}{1-\sigma} [\ln(s_g) - \ln(s_0)] \\ \Leftrightarrow (1-\sigma) \cdot [\ln(s_j) - \ln(s_0)] &= \delta_j - \sigma [\ln(s_g) - \ln(s_0)] \\ \Leftrightarrow (1-\sigma) \cdot \ln(s_j) + \sigma \ln(s_g) - \ln(s_0) &= \delta_j \end{aligned}$$

Finally, using the fact that $s_g = s_j/s_{j|g} \Leftrightarrow \ln(s_g) = \ln(s_j) - \ln(s_{j|g})$, we can get:

$$\Leftrightarrow \ln(s_j) - \sigma \ln(s_{j|g}) - \ln(s_0) - X_j\beta + \alpha p_j = \xi_j$$

2.5.3 Endogeneity issues

Simultaneity

Regardless of the model we use, we will have to deal with endogeneity issues regarding many dimensions, such as prices, observable characteristics or market shares.

In fact, in both models price was correlated to the unobservable term as firms make their decision based on demand which includes this term; moreover, in the nested logit case, we also have to deal with endogeneity in the within-share term. In general, anything that is correlated with the unobservable quality term ξ will have to be instrumented.

Measurement errors

Measurement errors in observed prices, characteristics or quantities may also create difficulties for the estimation procedure outlined above. Prices enter linearly in the estimation and are already endogenous so that these errors do not create too important biases. However, when these errors are present in market shares or quantities data then this issue becomes more important. Indeed, market share data is used to invert the model and get the unobservable term as a function of observed data. If shares include measurement errors, the non-linear inversion will accentuate the errors in a non-tractable way, creating a lot of issues in estimation. In practice, maximum likelihood will not have this kind of issue, but will make IV strategies harder.

Identification

For identification, we need to satisfy orthogonality conditions between ξ_j and the covariates (product characteristics x_j and price p_j). As discussed earlier, price will for sure be endogenous, so we need at least one instrument for this. Formally, we need a set of instruments w such that $E[\xi|x, w] = 0$.

In practice, the instruments for price will include “cost-shifters” (i.e. variables

that affect price without affecting demand).

2.5.4 Adding supply-side restrictions

It is also possible to get more orthogonality conditions by adding a supply-side to the demand system. Essentially, this means adding the firms' decisions about prices as moment conditions. To do this, we have to assume two things: (1) the firms' cost functions and (2) the competition model. Then, using the FOCs of the firms, one can recover extra moments to use in combination of the ξ moments to estimate the demand system by GMM.

Chapter 3

Mixed Logit (BLP)

For now we have seen three types of discrete-choice models and their applications to demand estimation for differentiated products: the simple logit model, the nested logit model and the multinomial probit model. The first two models, although quite useful and fairly simple to estimate were limited in three dimensions: they do not allow for random taste variation (2.2.2), they have restricted substitution patterns (2.2.4), they do not allow for correlation over time. Mixed logit models (which contains BLP) are highly flexible models that can deal with the previously mentioned issues.

Mixed logit models are a class of models that encompasses all types of models where the market shares are computed as integrals over simple logit functional form. We'll see in detail later what that means but intuitively, you should think of the model as everyone having her/his own logit model of demand, and aggregate demand would be computed by integrating over consumer attributes. In particular, IO economists are most interested in the BLP (for [?]) model and extensions of it like [?].

We'll first cover the basics of mixed logit and random coefficients, before talking in depth about estimation techniques as found in BLP and DFS.

3.1 Model

Generally, we write utility derived from consumption of a good j by consumer i as the function $u_{ij}(X_j, p_j, \xi_j, v_i, \theta)$, which is a function of product j observable characteristics (X_j of dimension $L - 1$), price (p_j), unobservable characteristics (ξ_j) and consumer characteristics (observable z_i and unobservables v_i), all entering the utility function through a vector of parameters θ .

As a simpler case, define utility as a linear function of those parameters such that:

$$U_{ij} = \sum_{l=1}^L x_{jl} \beta_{il} + \xi_j + \varepsilon_{ij}$$

$$\text{where } \beta_{il} = \lambda_l + \beta_l^o z_i + \beta_l^u v_i$$

Notice the main differences with logit models as we know them: first, price is not separated but included in observable characteristics of product j because of the second point, that β_i is a coefficient dependent on consumer characteristics. This means that different consumers will have different tastes in the same characteristics. For example, some person might be interested in the price of a phone while someone would disregard this and focus only on memory, another one on camera quality, etc.

We can rewrite the utility function by plugging in the definition of β :

$$\begin{aligned} U_{ij} &= \sum_{l=1}^L x_{jl} (\lambda_l + \beta_l^o z_i + \beta_l^u v_i) + \xi_j + \varepsilon_{ij} \\ &= \underbrace{\sum_{l=1}^L x_{jl} \lambda_l + \xi_j}_{\delta_j: \text{product mean utility}} + \sum_{l=1}^L x_{jl} \beta_l^o z_i + \sum_{l=1}^L x_{jl} \beta_l^u v_i + \varepsilon_{ij} \end{aligned}$$

There are five elements in this utility function:

- Observed ($x_j \lambda$) and unobserved (ξ_j) product quality.
- Observed ($x_j \beta^o z_i$) and unobserved ($x_j \beta^u v_i$) consumer-product interactions.
- A type-I EV iid error term (ε_{ij}).

The main features of the mixed logit model rely on the consumer-product interactions. Because they are not restricted to take the logit form, they will display more reasonable substitution patterns when aggregated. With these interactions, products will be close in terms of the characteristics of the consumers who buy the product. For example, consider an auto market following a price increase for the BMW 7 series (very expensive and luxurious car). In the logit model, this would create substitution to all other cars based on their market shares, meaning that a small Toyota Echo would benefit from the price increase. In the nested logit, provided we defined nests correctly, only cars within the same nest would be affected (luxury cars). Already, this makes more sense, but in the mixed logit model, we use parameters estimated from price variation to figure out what happened to the people who were buying BMW 7 series.

Aggregate vs. micro data

We have seen that the general form of the mixed logit model includes both observed and unobserved consumer characteristics. In general however, market data does not come with exact description of consumer characteristics for each interaction. At best, we might have aggregate consumer data (about a geographical region, a point in time, ...) but most often we cannot observe any characteristic.

When we do observe consumer characteristics, we can use them in z_i to estimate the model. A particularly influential paper using this type of data is the MicroBLP paper. When such data is unobserved, we work with just the v_i part of the model, as in the original BLP paper. For the rest of this section, we assume that we do not observe z_i .

3.1.1 Building the shares

Recall the utility function from the previous section (without z_i by assumption):

$$U_{ij} = \delta_j + \sum_{l=1}^L x_{jl} \beta_l^u v_i + \varepsilon_{ij}$$

Using the fact the error term is a type-I EV as in the simple logit, we can integrate it and get logit shares, but only conditional on v_i ! Formally, we get:

$$P_{ij}|v_i = \frac{\exp(\delta_j + x_j\beta v_i)}{\sum_k \exp(\delta_k + x_k\beta v_i)}$$

$$\Rightarrow P_{ij} = \int \frac{\exp(\delta_j + x_j\beta v_i)}{\sum_k \exp(\delta_k + x_k\beta v_i)} f(v_i) dv_i$$

Latent-class models

In the simplest case, consider that there are only two types of people, such that only two β_i are possible, say β_H and β_L . There are also two probabilities associated with the types. Thus the shares become:

$$P_{ij} = \frac{\exp(\delta_j + x_j\beta_H)}{\sum_k \exp(\delta_k + x_k\beta_H)} \cdot p_H + \frac{\exp(\delta_j + x_j\beta_L)}{\sum_k \exp(\delta_k + x_k\beta_L)} \cdot (1 - p_H)$$

3.2 BLP algorithm

Consider the situation where data is available only on the aggregate product-level, and no consumer data is observable. We will work out the estimation of demand following four steps:

1. Work out the aggregate shares conditional on both δ (mean utility) and β (taste variation).
2. Invert the shares to get ξ .
3. Estimate the model using the method of moments.
4. Repeat until convergence of all parameters.

To help understanding the details of the following steps, you should keep in mind a quick summary of what it to be done. As always, we are interested in the parameters of utility (that affect demand), thus λ and β^u (that we now call β only). These parameters are estimated using the interaction of the unobservable product characteristic ξ and adequate instruments. Until now, nothing should surprise you as we follow the same steps as described the GMM estimation strategy

for logit and nested-logit models. However, this time it is different since to get ξ , you will need the parameters you are looking for (this is the main difference of BLP). That is why you use starting values for those, and iterate until you find the parameters that are stable through estimation (this is a fixed point problem).

Step 1: Conditional aggregate shares

Recall that in the simple and nested logit models we studied, the probability P_{ij} that a consumer i buys product j was equal to the market share since they did not differ from the representative consumer in any way. This time, we now that two different consumers will have different probabilities of buying the product. Consequently, the market share is going to be the integral of consumers individual probabilities over their characteristics (here v_i since we do not observe any consumer characteristics).

Recall that we computed the market shares as a function of product characteristics:

$$s_j(\delta, \beta) = \int \frac{\exp(\delta_j + X_j v_i \beta^u)}{1 + \sum_k \exp(\delta_k + X_k v_i \beta^u)} f(v) dv$$

The issue with this integral is that it cannot be solved analytically (v_i is usually a multivariate normal distribution, which makes the market share a multi-dimensional integral); we can approximate it by taking the sample average over a set of ns draws in a simulation. This yields:

$$\hat{s}_j^{ns}(\delta, \beta) = \frac{1}{ns} \sum_i \frac{\exp(\delta_j + X_j v_i \beta^u)}{1 + \sum_k \exp(\delta_k + X_k v_i \beta^u)}$$

which is a function of parameters that we want to estimate (δ and β). We will see that δ will be computed, and β will be estimated by doing this step multiple times and finding the best one. This is why you should keep the simulated v_i for the whole exercise, else it is never going to converge.

Note that integration over the distribution of v_i is a different problem than ε_{ij} for which we can use the form of the type-I EV distribution to help us with analysis. In the case of v_i , we assume a multivariate normal distribution (across multiple consumer unobserved characteristics). If the distribution (not observations, but at least the pdf) is observed, then we can draw from the said distribution.

Moreover, notice that the use of a finite number of draws in the simulation will create a new source of errors within our estimation routine. Enough simulation draws should help tampering this issue, although finding the good number of draws is not an exact science, but more like a tradeoff between computation speed and errors. Overall, numerical evaluation of integrals is a particular topic that is deep enough to think about it carefully.

Step 2: BLP inversion

We now need to recover the product unobservable term ξ_j . In the same way as we did in the simple logit models, we already have the link between δ and ξ , but δ is "buried" in a nonlinear fashion into the the market shares: we need to invert the market shares to get delta, thus ξ , as a function of the shares (rather than the opposite). Doing that requires a special trick that is at the very core of[?] contribution.

Their trick is to use the fact that the following system:

$$\delta_j^k(\beta) = \delta_j^{k-1}(\beta) + \ln(s_j) - \ln(\hat{s}_j^{ns}(\delta^{k-1}, \beta))$$

is a contraction mapping. To see it, understand that s_j is the observed market share, the exponent k represents the iteration process. In other words, by iterating over this function, the δ values will converge to the true value, $\delta^*(\beta, s, \hat{s})$, where s, \hat{s} are respectively the vectors of observed and estimated shares conditional on β . Finally, we can write:

$$\xi(\beta, s, \hat{s}) = \delta^*(\beta, s, \hat{s}) - X_j \lambda$$

and use this form to construct the moments.

Step 3: Constructing the moments

Now that we have recovered ξ as a function of β (and λ through β since there is a way to write λ as a function of β), we can construct the moment conditions for demand estimation. To do this, we can go the OLS route if no component of x is

endogenous but most probably we will go the IV route, using w as the instrument matrix. The moment condition would therefore be:

$$E [\xi_j(\beta)w_j] = 0$$

As always in GMM, we want to select β such that the average analog to the moment equation is the closest possible to 0.

Step 4: Algorithm iteration

The first three steps were performed for a given β , thus we now need to find the best β , using a nonlinear search over β .

3.2.1 Identification

In order to fully identify the model, we need four sources of variation:

1. Choice set variation:
2. Product characteristics variation:
3. Consumer characteristics variation:
4. Functional form:

3.2.2 Helping estimation

Adding data

Adding a supply side

3.3 DFS algorithm

Chapter 4

Product Availability

4.1 Introduction

As we have seen in the previous chapters, many common models of demand estimation assume that the set of available choices (the “availability set”) is the same for all agents. While this assumption simplified computations by a lot, it may be a poor approximation of reality (think of retail environments). In fact, many consumers in many settings might encounter stockouts of products they intended to purchase. Moreover, availability variation can also come from assortment decisions (different stores displaying different products) or imperfect consideration (consumers not being aware of some alternatives). If availability variation arise exogenously, it could provide helpful sources of identification but in all cases, ignoring this variation will bias estimation.

There are three main sources of availability variation, from most to least aggregated:

- Assortment selection
- Stockout events
- Limited consumer information

This chapter will focus mainly on the first two, in relation to demand estimation (note that these topics might also be important on the supply-side).

4.2 Assortment variation

Understanding the effects of assortment on sales is mostly useful for two groups of agents: economists and policymakers who use demand estimates to evaluate welfare; and firms who want to explore these issues to come up with optimal choices.

Assortment variation is very prevalent in markets. For example, Tenn and Yun (2008) show that in grocery stores, 16% to 33% of sales happen in a time of intermediate availability. Bruno and Vilcassim (2008) show that in the confections market, the median availability of products is only 80%. If we consider time-series, we can see even more variability as Tenn and Yun explain by showing that only a third of products are available within a year and a third are discontinued every year.

Depending on the dataset, we can separate the problem of estimating demand with assortment variation in two cases: (1) assortment variation is unobserved and (2) assortment variation is observed.

4.2.1 Unobserved assortment variation

There are three main papers that discuss the challenges of estimating demand with unobserved assortment variation, using market-level data. The first one (Tenn, 2009) examines the sufficient assumptions for consistent estimates of own and cross-price elasticities in a linear demand model in product space. The second one (Tenn and Yun, 2008) explores the biases that arise due to assortment variation in logit types of model. Finally, Bruno and Vilcassim (2008) propose a method to account for this variation using additional assumptions.

Tenn (2009)

Tenn (2009) explores the potential issues of estimating linear demand using aggregate market-level data while ignoring assortment variation. The paper comes up with three conditions under which such a model would still yield consistent estimates:

- All prices within a market are the same, or in other words, all stores carry products at the same price.
- The fraction of stores that carry a product within a market is fixed across markets.
- Demand for each product is independent of the assortment.

Note how these assumptions are very strict and will be violated in many settings.

They show that the plim of estimated elasticities (again, when assortment is not accounted for) is given by:

$$\bar{\varepsilon}_{jkm} = E \left[\varepsilon_{jkm}^g | j, k \in g \right] \cdot \Pr [k \in g | j \in g]$$

where ε_{jkm}^g is the cross-price elasticity of product j with respect to k , given the assortment g (the object of interest). Intuitively, this means that $\bar{\varepsilon}_{jkm}$ is an “attenuated” (multiplied by a probability between 0 and 1) average of the elasticity for consumers who were presented assortments containing both j and k . In that regard, it will be biased towards 0! In the particular case of own-price elasticities however, this bias is absent since $\Pr [j \in g | j \in g] = 1$.

Tenn and Yun (2008)

Tenn and Yun (2008) consider a more complex model of demand in a standard logit. Utility derived from choosing good j in market m is given by the equation seen in chapter 2: $U_{ijm} = X_{jm}\beta - \alpha p_{jm} + \xi_{jm} + \varepsilon_{ijm}$, which yields the same type of elasticities:

$$\varepsilon_{jmr} = -\alpha p_{jm}(1 - s_{jmr}) \text{ for all } j \in J_{mr}$$

$$\varepsilon_{jkmr} = \alpha p_{km} s_{kmr} \text{ for all } j, k \in J_{mr} : j \neq k$$

Which at first do not seem biased, but in fact are because of (1) price coefficient bias and (2) elasticity estimates bias!

The price coefficient α can be biased in two ways. First, there might be an omitted variable bias, that arises with the failure to account for heterogenous retailers. In fact, a homogenous model such as the ones seen in chapter 2, will ignore the variation in availability of products (X_{jm} does not include a measure of availability, since it is not measured), thus introducing bias in the estimates. The second source

of bias in the price coefficients is the choice of instruments. Indeed, we usually use input prices or any “supply-side” variables to instrument for α , however, these variables might also have an effect on the assortment. For example, consider a sudden wholesale cost shock on a particular brand: the retailer might not want to purchase the good that period and change its assortment!

The elasticity estimates can also be biased in two ways: composition bias and weighting bias. As in the linear demand model, composition bias arises because of the fact that we average over all consumers, while some might not have had access to the same assortment. Moreover, bias in that case could go either way, depending on availability of the product (high availability means upward bias, low availability means downward bias). Weighting bias is the analog of the previous bias in aggregate, considering that markets as entities might have stores with different assortments.

Bruno and Vilcassim (2008)

Bruno and Vilcassim (2008) present an approach to correct the bias in a random-coefficients multinomial logit model such as the one discussed in chapter 3. Starting with a BLP-style utility where $U_{ijm} = X_{jm}\beta_i - \alpha_i p_{jm} + \xi_{jm} + \varepsilon_{ijm}$, we assume no observed individual variation so that we can rewrite utility as the sum of mean utility $\delta_{jm} = X_{jm}\beta - \alpha p_{jm} + \xi_{jm}$ and a consumer-specific term $\mu_{jm}(\theta_i) = X_{jm}\Sigma_\beta v_i - \sigma_\alpha \eta_i p_{jm}$.

Then, given an availability vector $a_m = (a_{1m}, \dots, a_{Jm})$ of indicator variables (if product j is available in market m or not), market shares can be written as:

$$\bar{S}_{jm}(a_m) = \int \frac{a_{jm} \cdot \exp(\delta_{jm} + \mu_{jm}(\theta_i))}{1 + \sum_{k=1}^J a_{km} \cdot \exp(\delta_{km} + \mu_{km}(\theta_i))} dG(\theta_i)$$

Then, given a probability mass function for a_m , denoted $\pi(a_m)$, the observed market shares are given by:

$$S_{jm} = \sum_{a_m} \bar{S}_{jm}(a_m) \pi(a_m)$$

To implement this, we need to add an outer loop in between the actual outer loop of BLP and the inner loop, in which we simulate the joint distribution of a_m to get the actual shares.

However, this method relies on two particularly “shaky” assumptions:

- Consumers’ tastes do not vary systematically with availability.
- Marginal distributions of availability are mutually independent.

4.2.2 Observed assortment variation

Including observed assortment variation looks like the solution to all our problems from the previous sections. However, Draganska, Mazzeo and Seim (2009) show that other challenges arise even when assortment is observed. Their method is analogous to the one described in Bruno and Vilcassim (2008), but does not relate on simulating availability since it is now observed for each market. However, if the unobservable quality of a good affects its presence in the assortment, the econometrician cannot identify δ_{jm} from ξ_{jm} , thus cannot use the BLP approach as in Bruno and Vilcassim. To go around that issue, they assume ξ_{jm} is simply uncorrelated with either price or assortment.

This issue is called assortment endogeneity and is treated in three more recent papers. The first one considers a simple reduced-form equation to recover availability of a product as a function of characteristics, and uses it as a control to relieve ξ_{jm} from its endogeneity. The second adds a reduced-form pricing equation. Finally, the third approach models a structural supply-side in order to solve for the equilibrium.

Iaria (2014)

Iaria (2014) deals with endogeneity of assortment by estimating a reduced form equation for assortment such that given product and market data, it is possible to recover if a product will be available or not. Note that this strategy, while not structural, is still valid since we only care about the demand side in this paper. Nonetheless, this approach relies on multiple regularity assumptions for identification.

Shah, Kumar and Zhao (2015)

Shah, Kumar and Zhao (2015) differs from Iaria (2014) in two main ways: first, they add a reduced-form pricing choice to control for endogeneity in that direction too and second, they estimate demand and the two reduced-form equations as a system, simultaneously.

Musalem (2015)

Finally, Musalem (2015) uses a more structural approach by modeling a supply-side directly in the model. Solving the model using Mathematical Programming with Equilibrium Constraints (MPEC). The supply side is a two-stage game in which firms choose the assortment first, then the prices. While this approach is the most structural, it requires assumptions on the equilibria played on both sides, and a game structure that approximates the real world.

4.3 Stockouts

When a product stocks out, the observed sales for this particular product or its substitutes does not correspond to the underlying demand anymore. By ignoring this phenomenon, demand estimates would be biased as in the assortment variation case. In fact, demand for stocked out products will be underestimated, which we call censoring bias while demand for the available substitutes will be overestimated, which we call forced substitution bias. Stockouts also have an effect on price coefficients if the events are correlated with price variation (stockout during a sale for example). If both types of bias arise in the same event, we cannot sign the bias anymore and estimation is confounded. However, exogenous stock-out events could provide interesting data about demand, especially if we look at diversion, etc.

There are two types of dataset used in this part of the literature: “periodic inventory” where we observe sales periodically, thus aggregating all sales within a period and “perpetual inventory” where we observe individual transactions. In the first type, stockouts are only partially observed since, within a period,

we cannot identify which sale created the stockout. In the second type though, stockouts are fully observed. The rest of this chapter describes the differences in treating these two types of datasets.

4.3.1 Partially observed stockouts

Both papers presented next develop demand estimation approaches for use with periodic store data, with inventory levels. This implies that stockouts will only be partially observed (who got the last unit of a product is unknown). Thus, the efforts should be directed towards estimating the timing of the stockout, or a distribution of sales of other products.

Conlon and Mortimer (2013)

The key insight in Conlon and Mortimer (2013) is that discrete-choice models of demand imply a distribution for the sales of all products given availability. In other words, it is possible to recover the sales using only the purchase probability implied by the demand model (total sales are thus a sufficient statistic). Here, the sales is the missing data.

Without going too deep in the details, the approach first recovers a probability distribution of the number of alternative sales before a stockout and integrates over it to get the expected sales of a particular alternative after the stockout (given the sales made before the stockout). This approach uses combinatorial computations, which makes evaluating the model very cumbersome when multiple stockouts arise. However, evaluating demand for a large number of products is possible since you only need information on the product of interest and the stockouts to get a result.

Musalem et al. (2010)

In Musalem et al. (2010), they do not consider the amount of sales as missing data but rather the sequence of sales. Using the fact that multiple sequences of sales could lead to the same end-of-period inventory, they express a likelihood function

conditional on the missing sequences and sum it over markets to recover the parameters. However, in this case, adding products is an issue since it will increase the possible sequences a lot, while adding stockouts will not affect computation.

4.3.2 Observed stockouts

4.4 Limited consumer information

Chapter 5

Entry Models

5.1 Introduction

5.1.1 Industry structure

5.1.2 Concentration measures

5.2 Sutton (1991)

The main question of the book from Sutton (1991) revolves around how do markets evolve to be less or more concentrated? Moreover, Sutton tries to explain why advertising is higher in concentrated industries. In order to explore the answers to these questions, he looks at different market structures and analyzes what makes these structure coherent with the data.

5.2.1 Perfect competition

Consider a market with perfect competition (price-taking assumption), exogenous sunk costs and free entry. The minimum efficient level of production (in the long run) is where $p = \min AC$.

This assumption means that, in the long run, firms will produce only the quantity satisfying this assumption, not more, not less. This implies that the number of firms in a market only depends on the size of the market, M . In particular, if you denote the quantity produced by firms at the optimum as q^* , you will have that the optimal number of firms n^* is given simply by:

$$n^* = M/q^*$$

As $M \rightarrow \infty$, we also have $n^* \rightarrow \infty$, thus the concentration ratio will tend to 0.

The issue with this simplistic model is that it might not hold in many settings where (1) competition might be imperfect or (2) sunk costs might be endogenous. As an example consider the increase of population in the United States since the creation of Pepsi and Coca-Cola. Even though the population has more than doubled, their respective market shares are still about 30% and 60%, which clearly disproves the previous model. For that reason he explores both directions.

5.2.2 Imperfect competition

Now suppose you allow for imperfect competition. For that purpose, we consider a two-stage game where in the first stage, firms choose to enter a market and incur a cost of A ; then in the second stage, firms play the market game. Moreover, assume that all M consumers have an income of 1 to spend exclusively on the good produced by the firms, such that in equilibrium:

$$p \cdot Q = M$$

Finally, suppose the marginal cost of production is $c > 0$.

Let N be the number of firms entering the market, each firm's profit will be denoted by $\Pi(N, M)$. Under free entry, we can get the equilibrium number of

firms, N^* , such that:

$$\Pi(N^*, M) > A \text{ and } \Pi(N^* + 1, M) < A$$

or in words, N^* is the number of firms such that any potential entrant would incur a loss by entering.

Moreover, it seems straightforward to say that for a given M (market size), $\Pi(\cdot)$ is a decreasing function of N (more competitors decrease profits if demand is fixed); and for a given N , $\Pi(\cdot)$ is an increasing function of M (more demand yields higher profits if we have the same number of competitors). Thus, as $M \rightarrow \infty$, we can say that $N^* \rightarrow \infty$ and concentration will still go to 0.

5.2.3 Endogenous sunk costs

Finally, consider a market with endogenous sunk costs, meaning that sunk costs will vary as a function of the current market structure. Two examples of this are advertising (being shown on the top of Google results pages is harder with more competitors) and R&D (having a better medication formula is also harder when competition is high). To introduce this intuition, consider a three-stage game with (1) firms enter the market, (2) they advertise their products and (3) they play the market game.

Under free entry, the equilibrium number of firms, denoted by $N^*(M)$, is determined as:

$$\Pi(N^*, M) > A(N^*) \text{ and } \Pi(N^* + 1, M) < A(N^* + 1)$$

5.3 Bresnahan and Reiss (1991)

In a realistic research setting, it is hard to observe strategic variables on a market level (prices, costs, advertising expenses, for all firms). This is the reason why Bresnahan and Reiss (1991) innovated on an estimation using mainly superficial market observables: the number of firms (but not the shares) and other general market characteristics (population, income, etc.). The main assumptions of this model rely on a static Sutton model with symmetric firms and free entry.

Behavioral model

As in the Sutton-type of models, we solve it backwards, assuming the fixed costs of entry are incurred as sunk in the second period, when firms are choosing production.

Demand in the market m is given by:

$$Q_m = d(Z_m, p) \cdot S(Y_m)$$

where $d(\cdot)$ represents the demand function of a representative consumer and $S(\cdot)$ is the total number of consumers that would buy the product. Note that this demand specification has constant returns to scale: doubling S will double Q . Finally, we define the inverse demand curve as $P(Q, Z, Y)$.

Therefore in the second-stage under Cournot competition, each firm solves:

$$\max_{q_i} \Pi_{N,m} \equiv P(q_i, q_{-i}, Z_m, Y_m) \cdot q_i - F_N - C(q_i)$$

where F_N is the endogenous sunk cost associated with N firms entering the market. Without loss of generality (we proved a similar result in the Cournot-Sutton context), assume that in equilibrium, quantities will be symmetric such that $q_i = q_j = q^*$ for all i, j . Because N is already “chosen” in the second stage, we can write individual equilibrium production as $q^* \equiv d(Z_m, p) \cdot \frac{S(Y_m)}{N}$. Then, using this, the profits of each firm is given by:

$$\Pi_{N,m} = P(q^*, q^*, Z_m, Y_m) \cdot q^* - F_N - C(q^*) = \left(P_N - \frac{C(q^*)}{q^*} \right) \cdot d(Z_m, P_N) \cdot \frac{S}{N} - F_N$$

Now, for both the average variable cost $\frac{C(q^*)}{q^*}$ and the fixed cost F_N , let's allow them to be additively separable in components, one of them being dependent on the firm only (respectively AVC and F) and the other component being endogenous on the number of firms (respectively b_N and B_N), such that the total profit of an individual firm, facing a market with N total firms, is defined as:

$$\Pi_{N,m} = (P_N - AVC(q^*) - b_N) \cdot d(Z_m, P_N) \cdot \frac{S}{N} - F - B_N$$

In the first stage, using the free entry assumption, we know that if N^* firms enter a market m , it must be that $\Pi_{N^*,m} > 0$ and $\Pi_{N^*+1,m} < 0$. Conversely, we can look

at the entry threshold s_N , which is the minimum additional demand, not already covered by existing $N - 1$ firms that is required in order for the N th firm to enter. This threshold is the value of $s_N \equiv S_N/N$ such that:

$$\Pi_{N,m} = 0 \Leftrightarrow \frac{S_N}{N} = \frac{F + B_N}{(P_N - AVC_N - b_N)d_N}$$

We can finally also look at the ratio of successive thresholds:

$$s_{N+1}/s_N = \frac{F + B_{N+1}}{F + B_N} \frac{(P_N - AVC_N - b_N)d_N}{(P_{N+1} - AVC_{N+1} - b_{N+1})d_{N+1}}$$

Thus, in a fully competitive market, a firm would enter each time its cost to enter could be recovered in the market, making s_{N+1}/s_N tend to 1, as $N \rightarrow \infty$.

Empirical strategy

The threshold equations derived above ask for a lot of observed variables, namely prices, costs (variable and sunk) and more. In reality, it is very difficult to observe all these at the same time for all firms in a market, thus we need a model that could allow for less information, which is the essence of [?].

Consider a reduced-form profit function given by:

$$\Pi_{N,m} = \underbrace{S(Y_m, \lambda)}_{\text{size}} \cdot \underbrace{V_N(Z_m, W_m, \alpha, \beta)}_{\text{variable profit}} - \underbrace{F_N(W_m, \gamma)}_{\text{fixed}} + \epsilon_m$$

$\bar{\Pi}_{N,m}$

which has an endogenous component $\bar{\Pi}_{N,m}$, and an error-term ϵ_m that is not dependent on N . This should raise some memories to anyone who covered the demand estimation part of the course. In fact, assuming a specification on the error term reveals the model as an ordered probit model, very similar to the logit or multinomial probit models. To see this, consider a market with N firms, this means that:

$$\bar{\Pi}_{N,m} + \epsilon_m > 0 \text{ and } \bar{\Pi}_{N+1,m} + \epsilon_m < 0$$

which is equivalent to:

$$\bar{\Pi}_{N,m} > \epsilon_m > \bar{\Pi}_{N+1,m}$$

Assuming ϵ_m is drawn from an iid normal distribution with mean 0 and variance σ^2 , we have that:

$$\Pr[N_m] = \begin{cases} \Phi(\bar{\Pi}_{N,m}) - \Phi(\bar{\Pi}_{N+1,m}) & \text{if } N_m > 1 \\ 1 - \Phi(\bar{\Pi}_{2,m}) & \text{if } N_m = 1 \end{cases}$$

Before going on to the estimation of the model, we need to look at the reduced-form of each component of the profit function:

- The market size includes a combination of market size characteristics Y_m , such as the population, the neighbouring population, growth, commuting possibilities, etc.
- The variable profit per capita is defined as: $V_N = \alpha_1 + X\beta - \sum_{n=2}^N \alpha_n$ where:
 - X is a vector of relevant economic variables such as demand characteristics for the product (Z) and cost-shifters (W).
 - $\alpha_n \geq 0$ is an intercept component such that each firm entering a market has a negative effect on profits.
- The endogenous fixed costs defined as: $F_N = \gamma_1 + \gamma_L w_L + \sum_{n=2}^N \gamma_n \cdot \gamma_n$

Finally, the estimation is performed using maximum likelihood. In fact, for each observation of the market, we use the probability function defined above. This gives us a likelihood function as a function of all observed variables described in the list above. Taking the log of it yields the log-likelihood to be maximized in order to estimate the parameters of interest.

5.4 Berry (1992)

The previous paper was lacking heterogeneity in the sense that its results applies for the case where all firms have the same costs, same continuation payoffs, etc. A step in the direction of allowing some differentiation was made by [?], where fixed costs can vary across firms. Thus, the behavioral model stays identical, but the profit of firm k can now be divided into a common component $v_m(N)$, incurred by all firms in the same way, and an idiosyncratic component $\phi_{m,k}$ applying only to firm k . In particular, we have:

$$\Pi_{N,m,k} = \underbrace{X_m\beta - \delta \ln(N) + \rho u_{m,0}}_{v_m(N)} + \underbrace{Z_k\alpha + \sqrt{1 - \rho^2} \cdot u_{m,k}}_{\phi_{m,k}}$$

Note that in this equation, the combination of the terms $\rho \cdot u_{m,0} + \sqrt{1 - \rho^2} \cdot u_{m,k}$ actually represents the error term, denoted $\epsilon_{m,k}$. In this setting, ρ represents the correlation between error terms across firms in a given market. Finally, [?] assumes that:

$$\epsilon_{m,k} \sim N(0, \Sigma)$$

where Σ is a matrix with all off-diagonal terms equal to ρ .

As in the demand estimation case, this error term is known to the players (the firms) but not to the econometrician. We are still in the static, full information game where all firms know everything about all other firms. However this time we are not in an ordered probit setting since the very structure of the problem is endogenous to the number of firms. To see that, recall the error of the previous model did not include any other subscript than m , while this one includes something about k which is intrinsically linked to the number of firms.

Contrary to the previous models, this one can display multiplicity of equilibria. In fact, although the equilibrium number of firms is unique, the exact firms that enter the market can be different. This implies that we lose information on which firms enter a market. For example, the model might predict that in a given market, two firms will ultimately enter, but you could have firms 1 and 2, firms 2 and 3 or firms 1 and 3. In order to control that issue, the author assumes that firms enter in the order of their profitability. Using that assumption, we can simplify the paper and compute the probability of observing N firms in the market as:

$$\begin{aligned} \Pr[n_m = N | Z_m] &= \Pr \left[\epsilon_{m,1}, \dots, \epsilon_{m,K_m} : \sum_{k=1}^{K_m} \mathbb{I}\{v_m(N, Z_m) > \phi_{m,k}\} = N \right] \\ &= \underbrace{\int \dots \int}_{K_m \text{ times}} \mathbb{I} \left\{ \sum_{k=1}^{K_m} \mathbb{I}\{v_m(N, Z_m) > \phi_{m,k}\} = N \right\} dF(\epsilon_{m,1}, \dots, \epsilon_{m,K_m}, \theta) \end{aligned}$$

which is a multi-dimensional integral (of K_m dimensions) and as we've seen in the previous two chapters, it is hard to compute. In order to circumvent this computational issue, we could use the same method as in BLP: an inner loop simulating the sample moment of the integral for each value of θ , and an outer loop solving for the value of θ that minimizes the square distance between the observed and simulated number of firms. This technique is called the Simulated Nonlinear Least Squares (or SNLS).

Formally, the outer loop finds

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{M} \sum_m (N_m - E [n_m | Z_m, \theta, \epsilon_{m,1}, \dots, \epsilon_{m,K_m}])^2$$

where the expectation term is approximated by its simulated sample average over random iid draws of $(\epsilon_{m,1}, \dots, \epsilon_{m,K_m})$. The simulated sample average is defined as:

$$E [n_m | Z_m, \theta, \epsilon_{m,1}, \dots, \epsilon_{m,K_m}] \approx \frac{1}{S} \sum_s n_m^{*,s}(\epsilon^s, \theta)$$

$$\text{where } n_m^{*,s}(\epsilon^s, \theta) \equiv \max_{0 \leq n \leq K_m} \left\{ n : \sum_{k=1}^{K_m} \mathbb{I}\{v_m(n, Z_m) > \phi_{m,n} | \epsilon^s\} \geq n \right\}$$

or in words, the maximum number of firms n , such that the number of firms entering the market with a positive profit is greater than n .

5.5 Two-period models

Two-period models have been introduced for the first time in entry models, but they have been used in several other settings since, for example in models of investment, contracting, etc. Typically, these models rely on:

- A first period where the agents choose a state variable that will determine the nature of the game in the second period.
- A second period where the game is played following the state of the world decided on the first.

The solution concept used to determine the outcome of the game is called "sub-game perfection" which means a Nash equilibrium is chosen at each step played in the game. We also call this equilibrium the Subgame Perfect Nash Equilibrium, or SPNE. In the simple case of a two-period game, the SPNE is quite easy to determine by backwards induction:

- You solve for the NE of the second period game for each potential game played following first period choices.
- Assuming the outcomes computed above are realized, you solve for the first period choice that is a NE.

Consider a typical two-period à la Stackelberg game where a firm chooses a level of capital ($K_1 \geq 0$) in the first period, and the second firm chooses its own level ($K_2 \geq 0$) in the second period. A firm that chooses a capital level of 0 is choosing to stay out of the market. A firm that chooses a level strictly greater than 0 enters the market, although firm 2 would have to pay a fixed cost of entry equal to F .

The profits of firm 1 are:

$$\pi_1(K_1, K_2) = K_1 \cdot (1 - K_1 - K_2)$$

while the profits of firm 2 are:

$$\pi_2(K_1, K_2) = \begin{cases} K_2 \cdot (1 - K_1 - K_2) - F & \text{if } K_2 > 0 \\ 0 & \text{if } K_2 = 0 \end{cases}$$

In order to solve the game, we proceed by backwards induction. In the second period, we consider the decision of firm 2 over the level K_2 . It will choose a level $K_2 > 0$ if and only if the profit associated is greater than 0, meaning that:

$$K_2(1 - K_1 - K_2) - F \geq 0 \Leftrightarrow K_2(1 - K_1 - K_2) \geq F$$

In that case, firm 2 will choose:

$$K_2^* = \arg \max_{K_2} K_2(1 - K_1 - K_2) - F$$

which is the solution to the following FOC:

$$1 - K_1 - 2K_2^* = 0 \Leftrightarrow K_2^* = \frac{1 - K_1}{2}$$

Plugging this solution into the participation condition above, we get that:

$$\frac{1 - K_1}{2} \left(1 - K_1 - \frac{1 - K_1}{2}\right) \geq F \Leftrightarrow \frac{(1 - K_1)^2}{4} \geq F \Leftrightarrow K_1 \leq 1 - 2\sqrt{F}$$

and thus the optimal level K_2^* is defined as:

$$K_2^* = \begin{cases} \frac{1 - K_1}{2} & \text{if } K_1 \leq 1 - 2\sqrt{F} \\ 0 & \text{if } K_1 > 1 - 2\sqrt{F} \end{cases}$$

Now that the second period is solved, we can go into the first period. Firm 1 will choose a level K_1 to maximize its profits, which are now known to be:

$$\pi_1(K_1) = \begin{cases} K_1 \cdot \left(1 - K_1 - \frac{1-K_1}{2}\right) & \text{if } K_1 \leq 1 - 2\sqrt{F} \\ K_1 \cdot (1 - K_1) & \text{if } K_1 > 1 - 2\sqrt{F} \end{cases}$$

From this perspective, it is clear that, for a given K_1 , allowing the other player to enter the market is worse than deterring entry. However, deterring entry can only be done for $K_1 > 1 - 2\sqrt{F}$. This creates a discontinuous profit function where deterring entry is always better but can be achieved only at some points that might yield a lower profit than allowing firm 2 to enter. The following graph represents this discrepancy nicely:

5.6 Entry models with multiple equilibria

In this section, we will reconsider static entry games (with two firms for simplicity) in order to introduce new ideas related to partial identification. In fact, as we have already discussed in this chapter, entry games have the potential to display multiple equilibria, which make straightforward estimation impossible. We have to move to a new way of estimation called “partial identification” of structural parameters. This concept simply means that instead of having a definite value for our structural parameters, we will recover multiple values, all coinciding with the model and the data. To recover these sets, we will need a new type of estimating equations called “moment inequalities”, which is a very recent topic in econometrics.

There are two ways to get to these moment inequalities: (1) using “structural” errors (as in Ciliberto and Tamer, 2009) or (2) using expectational or measurement errors (as in Pakes, Porter, Ho and Ishii, 2015).

5.6.1 Entry games with structural errors

Consider a simple two-firm entry model. Let $a_i \in \{0, 1\}$ denote the action of player $i = 1, 2$. The profits are given by:

$$\Pi_i(s) = \begin{cases} \beta' s - \delta a_{-i} + \varepsilon_i & \text{if } a_i = 1 \\ 0 & \text{if } a_i = 0 \end{cases}$$

where $\beta' s$ is the market profits (does not depend on the number of firms). We can see here that choices will be interdependent as profits are linked to the competitor's choice as well as the agent's. Moreover, note that there is an error term ε_i which is “structural”, meaning that both firms observe it, but the econometrician does not. In that sense, this model is a perfect information game.

Solving the game

For fixed values of the errors and parameters, the Nash equilibrium values a_1^* and a_2^* will satisfy best-response conditions, which are as follows:

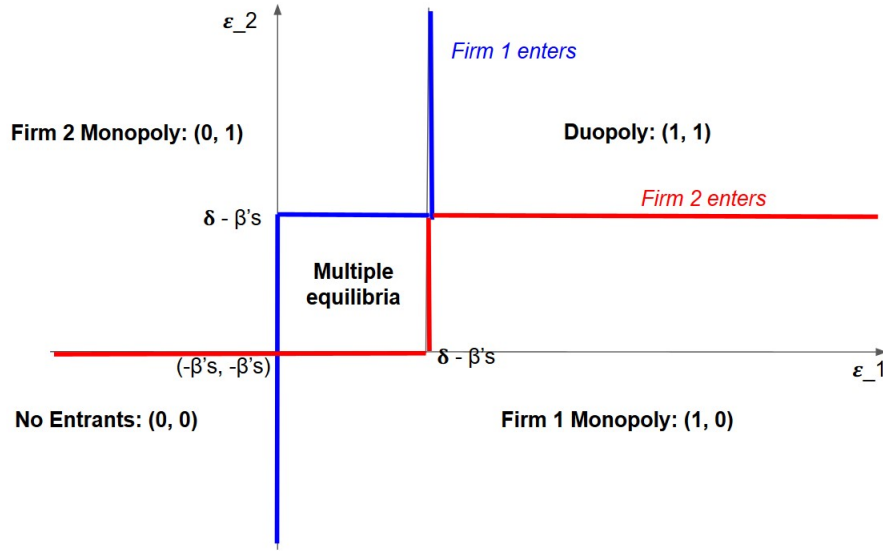
$$a_1^* = 1 \Leftrightarrow \Pi_1(a_2^*) \geq 0 \Leftrightarrow \beta' s - \delta a_2^* + \varepsilon_1 \geq 0$$

$$a_1^* = 0 \Leftrightarrow \Pi_1(a_2^*) \leq 0 \Leftrightarrow \beta' s - \delta a_2^* + \varepsilon_1 \leq 0$$

$$a_2^* = 1 \Leftrightarrow \Pi_2(a_1^*) \geq 0 \Leftrightarrow \beta' s - \delta a_1^* + \varepsilon_2 \geq 0$$

$$a_2^* = 0 \Leftrightarrow \Pi_2(a_1^*) \leq 0 \Leftrightarrow \beta' s - \delta a_1^* + \varepsilon_2 \leq 0$$

Using these, we can see that under some values of parameters, we might find multiple equilibria. Intuitively, think of a situation in which the market can hold only one firm, so that both $(1, 0)$ and $(0, 1)$ are possible equilibria. Alternatively, this picture shows this fact clearly:



Deriving the inequalities

Given this setup, we can derive inequalities for the probabilities of observing any state of the markets, using the distribution of the structural errors:

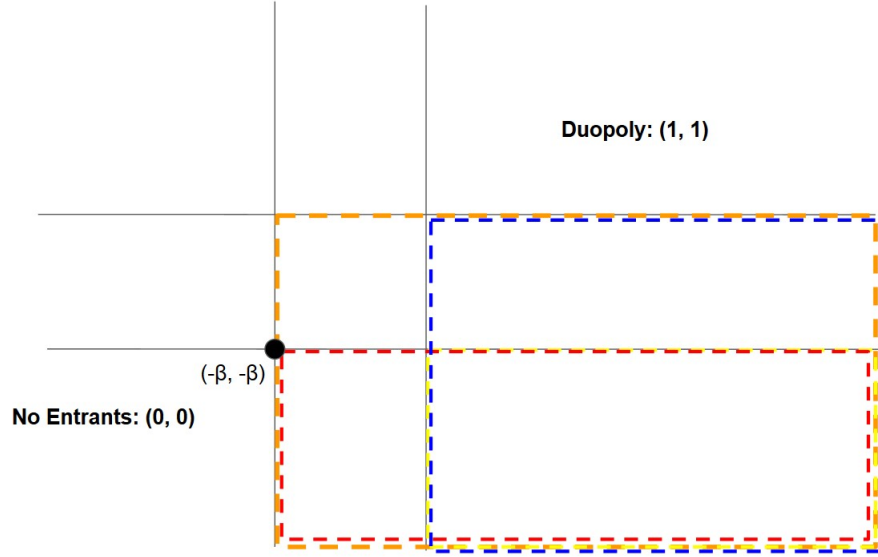
The probability of observing none of the two firms in the market is given by:

$$\Pr [a_1^* = a_2^* = 0 | s] = \Pr [\varepsilon_1 \leq -\beta' s \text{ and } \varepsilon_2 \leq -\beta' s] = F(-\beta' s)^2$$

For the two firms to be observed in the market:

$$\Pr [a_1^* = a_2^* = 1 | s] = \Pr [\varepsilon_1 \geq \delta - \beta' s \text{ and } \varepsilon_2 \geq \delta - \beta' s] = (1 - F(\delta - \beta' s))^2$$

For a monopoly of firm 1 to be observed, we have an issue because of the central square in the graph, where we do not know which firm is supposed to be in the market. We can however find an interval for that probability, using the following graph:



In fact, the probability of finding a firm 1 monopoly must be lower than firm 1 always being a monopoly (even when we see multiple equilibria), which is the orange area in the graph above:

$$\begin{aligned} \Pr[\varepsilon_1 \geq -\beta's \text{ and } \varepsilon_2 \leq \delta - \beta's] &\geq \Pr[a_1^* = 1, a_2^* = 0|s] \\ \Leftrightarrow [1 - F(-\beta's)] \cdot F(\delta - \beta's) &\geq \Pr[a_1^* = 1, a_2^* = 0|s] \end{aligned}$$

On the other side, if all multiple equilibria go to firm 2, we would observe at least a probability equal to the blue and red areas, minus the yellow:

$$\begin{aligned} \Pr[\varepsilon_1 \geq -\beta's \text{ and } \varepsilon_2 \leq -\beta's] + \Pr[\varepsilon_1 \geq \delta - \beta's \text{ and } \varepsilon_2 \leq \delta - \beta's] \\ - \Pr[\varepsilon_1 \geq \delta - \beta's \text{ and } \varepsilon_2 \leq -\beta's] &\leq \Pr[a_1^* = 1, a_2^* = 0|s] \\ \Leftrightarrow [1 - F(-\beta's)] \cdot F(-\beta's) + [1 - F(\delta - \beta's)] \cdot F(\delta - \beta's) \\ - [1 - F(\delta - \beta's)] \cdot F(-\beta's) &\leq \Pr[a_1^* = 1, a_2^* = 0|s] \end{aligned}$$

Obviously, the same computation can be done for the probability of observing only firm 2 in the market.

Now, define mutually exclusive outcome indicators as:

$$Y_1 = \mathbb{I}\{a_1 = 1, a_2 = 0\}$$

$$Y_2 = \mathbb{I}\{a_1 = 0, a_2 = 1\}$$

$$Y_3 = \mathbb{I}\{a_1 = 0, a_2 = 0\}$$

$$Y_4 = \mathbb{I}\{a_1 = 1, a_2 = 1\}$$

where a_i is the observed decision of the agent i . As such, we observe (or compute from observations) in the dataset the variables $Y_t = \{Y_{1t}, Y_{2t}, Y_{3t}, Y_{4t}\}$ for all $t = 1, \dots, T$. Using these variables, we can estimate the outcome probabilities $\hat{P}_{00}, \hat{P}_{01}, \hat{P}_{10}$ and \hat{P}_{11} (naively, without using any covariates), which we can then use in our inequalities from above to estimate the set of parameters (β, δ) .

Ciliberto and Tamer (2009)

5.6.2 Entry games with expectational errors

In contrast to the above, Pakes, Porter, Ho and Ishii (henceforth PPHI) derive the moment inequalities directly from the optimality conditions. Thus, the error term is not only unobserved to the econometrician but also to the firms.

Deriving the inequalities

In the two-firm entry game, if we observe actions (a_1^*, a_2^*) , then the Nash condition tells us that, given information sets I_1 and I_2 , we have:

$$E [\pi(a_1^*, a_2^*, z) | I_1] - E [\pi(a, a_2^*, z) | I_1] > 0, \text{ for all } a \neq a_1^*$$

$$E [\pi(a_1^*, a_2^*, z) | I_2] - E [\pi(a_1^*, a, z) | I_2] > 0, \text{ for all } a \neq a_2^*$$

where z is the vector of variables that have an effect on the profits (some might be unobserved to the agents). Because of that, PPHI suggests to parameterize the profit function as a function of actions and variables z , given a set of parameters θ to estimate. Formally,

$$\pi_i(a_1, a_2, z) = r_i(a_1, a_2, z; \theta)$$

which we can then use in the conditional moment inequalities from above to get:

$$E [r_1(a_1^*, a_2^*, z; \theta) - r_1(a, a_2^*, z; \theta)] > 0, \text{ for all } a \neq a_1^*$$

$$E [r_2(a_1^*, a_2^*, z; \theta) - r_2(a_1^*, a, z; \theta)] > 0, \text{ for all } a \neq a_2^*$$

Ho, Ho and Mortimer (2012)

Chapter 6

Moment Inequalities

This chapter provides a more rigorous introduction to estimation of models through the use of inequality restrictions, henceforth called moment inequalities. We have seen them in the previous chapter on entry, but they can be applied more generally to any type of games that would yield cumbersome computations using traditional methods, or for estimation when data are imperfect.

6.1 Framework

6.1.1 The agent's decision problem

Consider a situation with n decision makers indexed by i , having access to their own information set I_i when decisions are made and D_i the set of available decisions. The strategy played by agent i is a mapping $s_i : I \rightarrow D$ (from information to action), such that it generates the observed decisions d_i (which could be a vector).

The profit function of agent i is determined by his decision (d_i), the other agents' decisions (d_{-i}) and other environment variables y_i . At the time of the decision, the agent has expectations over what happens in the game ($\pi(\cdot)$, s_i , I_i and Y_i); they are denoted by $\mathcal{E}[\cdot]$, which is not the same operator as the typical expectation

operator.

Best-response condition (Nash)

If d_i is the observed decision of player i , we assume:

$$\sup_{d \in \mathcal{D}_i} \mathcal{E} [\pi(d, d_{-i}, y_i) | I_i] \leq \mathcal{E} [\pi(d_i, d_{-i}, y_i) | I_i] \text{ for all } i = 1, \dots, n$$

Quite obviously, we can see this assumption as an assumption for “rationality”, meaning that at the time of the decision, the agent chose the best option. In single agent problems, this comes directly from optimization behavior, while in games, it is only a necessary condition for a Bayes-Nash equilibrium to be played, but it does not rule out multiple equilibria, or restrict the selection between equilibria.

Counterfactual condition

In order for the agents to ensure optimal behavior, they need to evaluate the alternative decisions in their counterfactual environment. Thus we need to define what happens to d_{-i} and y_i following the decision of agent i . Note that in single-agent problems and simultaneous games, the counterfactual is assumed away using a conditional independence assumption.

In other cases, we assume that $y_i = y(z_i, d, d_{-i})$ and that the distribution of (d_i, z_i) conditional on I_i and d do not depend on d . In words, this assumption means that environment variables y_i depend only on variables z_i and decisions by the agents (which are all exogenous conditional on I_i and d).

Using this, we define the “differential profit” as:

$$\Delta\pi(d, d', d_{-i}, z_i) = \pi(d, d_{-i}, y(z_i, d, d_{-i})) - \pi(d', d_{-i}, y(z_i, d', d_{-i}))$$

as the difference in profits between two decisions.

Finally, we rewrite the first condition as:

$$\mathcal{E} [\Delta\pi(d_i, d, d_{-i}, z_i) | I_i] \geq 0 \text{ for all } i = 1, \dots, n$$

This might seem like the inequality to use in estimation, however, recall that the expectation is only the agent's so we need to recover empirical analogues of these in order to use them.

6.1.2 Observables and disturbances

We assume that the econometrician has a parametric function, denoted $r(\cdot)$, that approximates $\pi(\cdot)$ given arguments d_i , d_{-i} , observable variables of z_i , denoted z_i^o and unknown parameters θ to estimate.

Using that function, we can approximate the differential profit with $\Delta r(d, d', d_{-i}, z_i^o, \theta)$. From there, define two types of errors:

$$\begin{aligned} v_{2,i,d,d'} &= \mathcal{E} [\Delta \pi(d, d', d_{-i}, z_i) | I_i] - \mathcal{E} [\Delta r(d, d', d_{-i}, z_i^o, \theta) | I_i] \\ v_{1,i,d,d'} &= v_{1,i,d,d'}^\pi - v_{1,i,d,d'}^r \\ v_{1,i,d,d'}^\pi &= \Delta \pi(d_i, d, d_{-i}, z_i) - \mathcal{E} [\Delta \pi(d_i, d, d_{-i}, z_i) | I_i] \\ v_{1,i,d,d'}^r &= \Delta r(d, d', d_{-i}, z_i^o, \theta) - \mathcal{E} [\Delta r(d, d', d_{-i}, z_i^o, \theta) | I_i] \end{aligned}$$

that we refer to in general as $v_{2,i}$ and $v_{1,i}$ (composed of $v_{1,i}^\pi$ and $v_{1,i}^r$). Note that the first error, while not observed by the econometrician, is a part of the information set I_i of the agent (he “knows” $v_{2,i}$). The second error is not observed by either the agent nor the econometrician.

6.1.3 Moment inequalities

6.2 Applications

Chapter 7

Single-agent Discrete Dynamic Programming

7.1 Definition

A discrete dynamic programming is a maximization problem of an objective function, over an infinite horizon, with discounting. More formally, the objective function takes the following form:

$$E \left[\sum_{t=0}^{\infty} \beta^t \cdot r(s_t, a_t) \right]$$

where s_t is the state variable, a_t is an action, β is the discount factor and $r(\cdot)$ is the “reward” function for taking action a_t in the state s_t . While the action is deliberately taken by the optimizing agent, the state is pinned down by previous actions and states. We say that at each point in time, the state summarizes the information about the world. Hence, s_t not only contains information about the present time period but also all previous time periods. This means that each pair (s_t, a_t) will determine the potential state s_{t+1} ; we denote the probability of moving to state s_{t+1} given taking action a_t in state s_t as $Q(s_t, a_t, s_{t+1})$. In other words, the actions that the agent takes influence not only current rewards but also the future time path of the state, and thus future rewards as well. The essence of dynamic programming problems is solving this trade-off between current and

future rewards at each point in time.

Policies

The most fruitful way to think about solutions to DDP problems is to compare what we call policies. A policy is a mapping from past actions and states to current action. In other words, it is a sort of code of conduct for the agent. Recall that the current state holds sufficient information about past actions and past states, therefore, a policy can also be represented as a mapping from current state to current action. This type of policies is called stationary Markov policies. A stationary Markov policy (henceforth just policy) is denoted as $\sigma(\cdot)$ which maps actions to states such that $a_t = \sigma(s_t)$. It is known that, for any arbitrary policy, there exists a stationary Markov policy that dominates it at least weakly: this is why solving the previous problem can be reduced to finding the best policy, only among Markov stationary policies.

Formal definition

A discrete dynamic program consists of:

- A finite set of states $S = \{1, \dots, n\}$.
- A finite set of feasible actions $A(s)$ defined for each $s \in S$, and thus a set of feasible state-action pairs: $SA = \{(s, a) : s \in S, a \in A(s)\}$.
- A reward function $r : SA \rightarrow \mathbb{R}$.
- A transition probability function $Q : SA \rightarrow \Delta S$, where ΔS is the set of probability distributions over S .
- A discount factor $\beta \in [0, 1)$.

Over this definition, we provide other definitions to help us solve the problem. First, we define a policy as a function $\sigma : S \rightarrow A$. Then, in the set of policies are the feasible policies such that $\sigma(s) \in A(s)$ for all $s \in S$. This subset is denoted Σ .

Now that we have a structure, let's remind ourselves how the problem works. If an agent uses a policy $\sigma \in \Sigma$, then he gets the current reward $r(s_t, \sigma(s_t))$ at time t and the probability that next period's state is s' is given by $\Pr[s_{t+1} = s'] = Q(s_t, \sigma(s_t), s')$.

Then, for each feasible policy $\sigma \in \Sigma$, define:

- $r_\sigma = r(s, \sigma(s))$, the vector of possible rewards at each state, for a given policy.
- $Q_\sigma = Q(s, \sigma(s), s')$, the state-transition matrix from s to s' , also called the stochastic matrix on S .

Using these definitions, note that a row of Q_σ contains the probabilities of all future states given the state of the row happened. For example, the s -th row would be a row vector of probabilities of each state s_i happening in the next period, which would sum to 1. This is interesting because, coupled with r_σ , we get that:

$$(Q_\sigma \cdot r_\sigma)(s) = \mathbb{E}[r(s_t, \sigma(s_t)) | s_{t-1} = s]$$

where $(Q_\sigma \cdot r_\sigma)(s)$ denotes the s -th row, of the column vector $Q_\sigma \cdot r_\sigma$ and assuming $s_t \sim Q_\sigma$.

Value and Optimality

Let $v_\sigma(s)$ denote the discounted sum of expected rewards following policy σ , given an initial state of s , or formally:

$$v_\sigma(s) = \sum_{t=0}^{\infty} \beta^t (Q_\sigma \cdot r_\sigma)(s)$$

This function is called the policy value function for σ .

A policy $\sigma \in \Sigma$ is optimal if, for all $s \in S$, we have that:

$$\sigma \in \arg \max_{\sigma \in \Sigma} v_\sigma(s)$$

The optimal value function, or simply value function, is, as its name suggests, the value function of the previous problem, or formally, $v^* : S \rightarrow \mathbb{R}$ such that:

$$v^*(s) = \max_{\sigma \in \Sigma} v_\sigma(s)$$

Finally, given any value function $w : S \rightarrow \mathbb{R}$, a feasible policy $\sigma \in \Sigma$ is called w -greedy if

$$\sigma(s) \in \arg \max_{a \in A(s)} r(s, a) + \beta \sum_{s' \in S} w(s') Q(s, a, s') \quad \text{for all } s \in S$$

This means that, given an expected future path of rewards determined by w , the w -greedy policy σ is a policy that yields the best current action to take, for each state. Hence, it follows that if we knew the optimal value function v^* , we would be looking for policies that are v^* -greedy, or that would suggest the best current action to take given following the optimal path in the future.

Operators

Before going further into solving the problem, we define two useful operators.

The Bellman operator, $T : \mathbb{R}^S \rightarrow \mathbb{R}^S$ is defined by:

$$(Tv)(s) = \max_{a \in A(s)} r(s, a) + \beta \sum_{s' \in S} v(s') Q(s, a, s') \quad \text{for all } s \in S$$

This operator transforms any value function by changing its current action to the optimal one in any state. In other words, if your current value function is bad, the Bellman operator will make it slightly better by altering the current action to its optimal value, thus increasing the level of the overall value function. In that sense, the Bellman operator is monotonic, such that if $v \leq w$, then $Tv \leq Tw$. Moreover, we can show that the Bellman operator is also a contraction mapping, or formally that: $\|Tv - Tw\| \leq \beta \|v - w\|$. Using these two facts, we will be able to design an iterative algorithm to get to the solution of the problem.

The other operator is ...

Bellman equation and Optimality

We know that the solution of the DDP problem is to find a value function such that:

$$v(s) = \max_{a \in A(s)} r(s, a) + \beta \sum_{s' \in S} v(s') Q(s, a, s') \quad \text{for all } s \in S$$

which can be written using the Bellman operator as:

$$v(s) = Tv(s) \quad \text{for all } s \in S$$

In this form, it is now clear that the essence of the problem is to find the fixed-point of the operator T .

Moreover, using the fact that we know T as a contraction mapping of modulus β , we get a natural way to get the fixed-point of T by applying it infinitely many times to any initial value function v^0 . Formally,

$$v^* = \lim_{k \rightarrow \infty} T^k v^0$$

However, in real applications, there is no such thing as infinite operations, thus we will need a finite approximation of this procedure.

Value Function Iteration

The value function iteration procedure is a finite approximation of the infinite Bellman operation we discussed just above. In fact, it relies on applying the Bellman operator as many times as needed to get a value function that is stable enough (close to the fixed point). The procedure goes as follows:

1. Choose any $v^0 \in \mathbb{R}^n$, and specify an error $\varepsilon > 0$ (close to 0). This is the initial step, $k = 0$.
2. Compute $v^{k+1} = Tv^k$, meaning that you need to solve the maximization problem defined as:

$$v^{k+1}(s) = \max_{a \in A(s)} r(s, a) + \beta \sum_{s' \in S} v^k(s') Q(s, a, s')$$

for each state $s \in S$.

3. Evaluate the continuation condition:

$$\sup_s \|v^{k+1}(s) - v^k(s)\| < \left\lceil \frac{(1 - \beta)}{2\beta} \right\rceil \varepsilon$$

If it is true, then go to the next step, otherwise go back to step 2 and set $k = k + 1$.

4. The optimal value function computed is v^{k+1} and the optimal policy function is σ , a v^{k+1} -greedy policy.

By performing this procedure, you get a ε -approximation of the optimal value function.

7.2 Rust (1987)

Harold Zurcher is the manager of a bus depot in Madison, WI. Each month, for each bus, he must decide if the bus should continue with the current engine or if it should be replaced. If the bus continues without replacing the engine, then it pays a constant marginal cost for each additional mile. If the bus gets replaced, then Harold pays a fixed replacement cost. The problem that Harold faces is to minimize long run expected cost.

This is a clear discrete state dynamic programming problem as defined above, and as such, should have been fairly easy to estimate at the time the paper was written. However, Rust (1987) is about estimating the parameters of the model, using the actual data from the decisions of Harold Zurcher. In that sense, solving the model is really the first step (the inner loop) of a more intricate process to recover the parameters of the model (fixed costs, marginal costs, etc.).

In this chapter, we will go over both the solution to the model (given a set of parameters) and the estimation of it. We draw on the problem set 4 given by Prof. Julie Mortimer.

7.2.1 Setup

Optimal stopping problem as a DDP

First we setup the model, using the formal definition of a DDP given in the first section. The model is set for a single bus (recall the assumption that there are no correlations between buses) for time periods $t = 1, 2, \dots, +\infty$ representing weeks (for each time HZ can take a decision on the bus).

The state variable s will be the mileage on the bus. Hence the set of states S should be the set \mathbb{R}_+ of all possible mileages.

In each state, there are two feasible actions, which are the same for any state: “not replace” (a^0) or “replace” (a^1). The feasible action set is constant for all states: $A(s) = \{a^0, a^1\} = A$. Thus, the feasible state-action pairs set is: $SA = \{(s, a) : s \in S, a \in A\}$.

The reward function is the profit function for each feasible state-action pair. For any state-action (s, a) where $a = a^1$, the profit will be a negative fixed replacement cost RC . For state-action pairs in which $a = a^0$, then the profit will be a negative cost depending on the mileage (the state s). Formally, we get:

$$r(s, a) = \begin{cases} -RC & \text{if } a = 1 \\ -c(s, \theta) & \text{if } a = 0 \end{cases}$$

The transition probability function yields a probability distribution over all states for each feasible state-action pair. To stay very general, we write:

$$s_{t+1}|s, a \begin{cases} = 0 & \text{if } a = 1 \\ \sim G(\cdot|s) & \text{if } a = 0 \end{cases}$$

where G is unknown up to parameters.

Finally, the discount factor is denoted as β .

Therefore, Harold Zurcher chooses the infinite sequence of actions $\{a_t\}_{t=1,\dots}$ to maximize the objective function:

$$\max_{\{a_t\}_{t=1,\dots}} \mathbb{E} \left[\sum_{t=1}^{\infty} \beta^t \cdot r(s_t, a_t; \theta) \right]$$

Define the (optimal) value function as:

$$V(s_t) = \max_{\{a_t\}_{t=1,\dots}} \mathbb{E} \left[\sum_{\tau=t}^{\infty} \beta^{\tau-t} \cdot r(s_\tau, a_\tau; \theta) | s_t \right]$$

And finally we write the Bellman equation:

$$a_t = \arg \max_a r(s_t, a; \theta) + \beta \mathbb{E} [V(s') | s_t, a]$$

Solving the model requires knowledge of everything described until now, then performing a value function iteration as described in the first section. However, there are parameters in the model that are not observed in the data, hence the need to estimate them. To estimate is different than to solve, thus we will need to add crucial elements and assumptions to the model.

Adding structural errors

In particular, the main addition required is a structural error in the reward function. This structural error will allow for “positive likelihood” given our imperfect model (because the econometrician cannot recover the function perfectly). Recall that a structural error is observed by the agent at the time of the decision, but is never observed by the econometrician.

We get that the “new” reward function is thus given by:

$$r(s, a) = \begin{cases} -RC + \varepsilon_1 & \text{if } a = 1 \\ -c \cdot s + \varepsilon_0 & \text{if } a = 0 \end{cases}$$

We are left with three types of parameters to estimate that are summarized in θ :

- Parameters of the $c(\cdot)$ function.
- Value of replacement cost RC .
- Parameters of the mileage transition function $G(\cdot|x)$.

Note that both the discount factor β and the distribution of structural errors are not to be estimated. In fact, most dynamic-discrete choice settings are shown to be nonparametrically underidentified. Intuitively, β is difficult to disentangle from RC .

7.2.2 Econometric model

We observe actions and mileages for $t = 1, \dots, T$ and 62 buses. We treat all buses as independent draws of the same distribution (as if all buses were equivalent).

The likelihood function for a bus, over T periods is given by:

$$\begin{aligned}
L(\theta) &= f(s_1, \dots, s_T, a_1, \dots, a_T | s_0, a_0, \theta) \\
&= \prod_{t=1}^T \Pr [s_t, a_t | s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1}, \theta] \\
&= \prod_{t=1}^T \Pr [s_t, a_t | s_{t-1}, a_{t-1}, \theta] \\
&= \prod_{t=1}^T \Pr [a_t | s_t, \theta] \cdot \Pr [s_t | s_{t-1}, a_{t-1}, \theta_3]
\end{aligned}$$

where θ_3 is the vector of the parameters that govern the distribution of mileages.

First of all, from the second to the third lines, we use the stationary Markov policy property (the best policies are Markov stationary, meaning that they only require information about the previous period).

Second, from the third to last equations, we use a new assumption called conditional independence. Formally, this assumption implies that (1) conditional on the current state, actions are independent over time (actions depend only on previous actions through the state variable) and that (2) conditional on the previous state and previous action, the current state is independent of parameters from the cost function or the structural errors.

From the likelihood derived above, in order to estimate the model, we first need to recover each element in the final equation:

1. Recover θ_3 to get $\Pr [s_t | s_{t-1}, a_{t-1}, \theta_3]$.
2. Estimate θ using $\Pr [a_t | s_t, \theta]$.

Markov transition probabilities (θ_3)

Instead of assuming a functional form for the mileage distribution $G(\cdot)$, Rust uses a simplifying trick where we assume a discrete distribution for the incremental

mileage $\Delta s_t = s_t - s_{t-1}$ such that:

$$\Delta s_t = \begin{cases} [0, 5000) & \text{w/ probability } p \\ [5000, 10000) & \text{w/ probability } q \\ [10000, \infty) & \text{w/ probability } 1 - p - q \end{cases}$$

so that $\theta_3 = \{p, q\}$. Note that this step does not depend on the rest of the model and can be estimated right away, nonparametrically.

Dynamic Logit

For the other part of the likelihood function, we use an additional assumption on ε_{1t} and ε_{0t} such that they are iid draws from two independent Type-I Extreme Value distributions (with the standard normalizations). This is why we call this part a dynamic logit.

Expanding the probabilities:

$$\begin{aligned} \Pr[a_t = 1|s_t, \theta] &= \Pr[r(s_t, 1; \theta) + \beta E[V(s')|s_t, 1] > r(s_t, 0; \theta) + \beta E[V(s')|s_t, 0]] \\ &= \Pr[-RC + \varepsilon_{1t} + \beta V(0) > -c \cdot s_t + \varepsilon_{0t} + \beta E[V(s')|s_t]] \\ &= \Pr[\varepsilon_{1t} - \varepsilon_{0t} > RC - \beta V(0) - c \cdot s_t + \beta E[V(s')|s_t]] \\ &= \frac{\exp(-RC + \beta V(0))}{\exp(-RC + \beta V(0)) + \exp(-c \cdot s_t + \beta E[V(s')|s_t])} \end{aligned}$$

as in the logit model! Obviously, this result applies also to the opposite action:

$$\Pr[a_t = 0|s_t, \theta] = \frac{\exp(-c \cdot s_t + \beta E[V(s')|s_t])}{\exp(-RC + \beta V(0)) + \exp(-c \cdot s_t + \beta E[V(s')|s_t])}$$

The main issue with this probability is that $E[V(s')|s_t]$ does not have a closed form (it is precisely what we want to solve in a DDP). This is why the estimation method will rely on compute this element in an inner loop, for every guess of parameters (outer loop). The next section explain this procedure.

7.2.3 Estimation Method

As mentioned in the previous section, the estimation procedure has to rely on two loops, because the expected value function is not known up to parameters θ . Thus, the first loop (the outer loop) will try to find the optimal θ given the value function computed in the second loop (the inner loop) which solves the model for a given θ .

Computational details

The goal of the second step is to compute the value function for a given value of θ . A clever and computationally convenient feature in the Rust paper is to perform value function iteration not on the value function per se, but on the expected value function denoted $EV(s, a)$ defined as:

$$EV(s, a) \equiv E[V(s')|s, a]$$

We can develop the expectation to get:

$$EV(s, a) = \int_{s'} \log \left(\sum_{a' \in \{0,1\}} \exp(r(s', a'; \theta) + \beta EV(s', a')) \right) \cdot \Pr[s'|s, a]$$

and now we see the parallel between this equation and the value function before the value function iteration of the first section. The only difference is that we do not need to “solve” for the best current action since we can use the logit inclusive value to recover the utility of the best choice! This yields a huge improvement in computational complexity.

Nested Fixed-point Algorithm

What follows after is the analog to the value function iteration. Let τ index the iterations, such that $EV^\tau(s, a)$ is the expected value approximation after τ iterations.

1. Choose any EV^0 and specify a tolerance $\eta > 0$ but very small. Usually, we choose $EV^0 = 0$.

2. Use $EV^{\tau-1}$ to compute:

$$\begin{aligned}
EV^\tau(s, a) &= \int_{s'} \log \left(\sum_{a' \in \{0,1\}} \exp(r(s', a'; \theta) + \beta EV^{\tau-1}) \right) \cdot \Pr[s'|s, a] \\
&= p \cdot \int_s^{s+5000} \log \left(\sum_{a' \in \{0,1\}} \exp(r(s', a'; \theta) + \beta EV^{\tau-1}) \right) ds' \\
&\quad + q \cdot \int_{s+5000}^{s+10000} \log \left(\sum_{a' \in \{0,1\}} \exp(r(s', a'; \theta) + \beta EV^{\tau-1}) \right) ds' \\
&\quad + (1-p-q) \cdot \int_{s+10000}^{\infty} \log \left(\sum_{a' \in \{0,1\}} \exp(r(s', a'; \theta) + \beta EV^{\tau-1}) \right) ds'
\end{aligned}$$

for each state and action (s, a) .

3. Evaluate the continuation condition:

$$\sup_{(s,a)} \|EV^\tau(s, a) - v^{\tau-1}(s, a)\| < \eta$$

If it is true, then we're done and EV^τ will serve as our best approximation (given θ), otherwise go back to step 2 and increment τ .

Using the newly computed function $EV(s, a)$, we can derive the probabilities that were unknown previously and come up with a value for the likelihood function. The outer loop will take that value and suggest another θ to compute a new approximation of the EV function, new probabilities and a new likelihood function value, until convergence to θ_0 : the true value of structural parameters.

7.3 Hotz-Miller approach

Chapter 8

Retailing and Inventories

8.1 Aguirregabiria (1999)

8.1.1 Summary

Background

There is substantial evidence of price dispersion and staggering (prices are not flexible), most commonly explained by menu costs (changing prices is costly) and not perfect correlation in demand (different demand shocks across firms). Confirmed by theory of (S, s) rule (= target price and adjustment bands with adjustment only outside the bands). But what to make of evidence that retail firms might not use (S, s) rules, or that prices go down sometimes (with inflation, prices should always go up)? Most explanations rely on “consumer-side” phenomena, not on inventories (supply-side)!

Maybe a model of (S, s) rule with inventories would explain some variation!

Model

Firms maximize profits as a function of expected sales, order costs, inventory costs and menu costs (all three are lump-sum!) = (S, s) rule is still optimal!

Without menu costs, inventories are perfectly negatively correlated with markups: giving evidence of price markdowns when inventory is high (right after an order) but no price staggering.

With menu costs, you get price staggering as well as markdowns.

Data

Monthly information on prices, sales, orders and inventories for every item = balanced panel data. Price data is monthly averages (for regular price and sales price).

Empirical evidence

Naive analysis and reduced form regressions show infrequent price changes, prevalence of sales promotions, infrequent orders, negative correlations between inventories and prices.

Using a Hotz-Miller approach, authors find the same results: data is consistent with high lump-sum cost of ordering and menu costs.

Counterfactuals

Removing these lump sum costs actually leads to more variability = they explain a lot of price staggering and price reductions! Demand might explain the rest.

8.1.2 Discussion

8.2 Hendel and Nevo (2002)

8.2.1 Summary

Background

When goods are storable, price variation might create bias in demand estimation. In fact, if prices decrease (during sales for example), then consumers might stockpile: creates large demand increase in the short run which, measured in long-run (when prices have gone back up) will inflate own-price elasticities. Cross-price elasticities be ambiguous. Because these elasticities are so important in industry analysis, you have to get them right!

This is the analog analysis to the Aguirregabiria (1999) paper where he studies price variation and inventories (= the supply side).

Data

Scanner data on store-level and household-level consumption in the detergent market. Hendel and Nevo observe price, quantities, some measure of advertising, sales, inventories (both at the store and household levels).

Reduced form analysis shows that duration between purchases has a positive effect on quantity purchased (= household do hold inventories); storage costs are negatively correlated with buying on sale (= not being able to store leads to less storage?); households buy more on sale, even when they do hold inventories already (= stockpiling).

Model

Do not model quantity purchases but only sizes! Dynamic discrete choice model with utility of purchase and inventory costs.