

SEMI-NONPARAMETRIC ESTIMATION OF CONSUMER SEARCH COSTS *

José Luis Moraga-González [†]
Zsolt Sándor [‡]
Matthijs R. Wildenbeest [§]

August 2011

Abstract

This paper studies the estimation of the distribution of the costs of non-sequential search. We show that the search cost distribution is identified by combining data from multiple markets with common search technology but varying consumer valuations, firms' costs, and numbers of competitors. To exploit such data optimally, we provide a new method based on semi-nonparametric (SNP) estimation. A Monte Carlo study shows that the method works well in relatively small samples. We apply our method to a data set of online prices for memory chips and find that the search cost density is essentially bimodal such that a large fraction of consumers searches very little, while a smaller fraction of consumers searches a relatively large number of stores.

Keywords: consumer search, oligopoly, search costs, semi-nonparametric estimation

JEL Classification: C14, D43, D83, L13

*We are indebted to Thierry Magnac, the Editor, two anonymous referees, Vladimir Karamychev, Allard van der Made, Paulo K. Monteiro, Martin Pesendorfer, Michael Rauh, and Matthew Shum for their useful comments. The paper has also benefited from presentations at Universidad Carlos III de Madrid, University of Chicago, CORE, University of Essex, Tinbergen Institute, the ESEM Meetings 2006 (Vienna), the EEA Meetings 2007 (Budapest), and the IIOC 2009 (Boston). Financial support from the Netherlands Organization for Scientific Research (NWO) and from Marie Curie Excellence Grant MEXT-CT-2006-042471 is gratefully acknowledged. Earlier drafts were circulated under the title "Nonparametric Estimation of the Costs of Non-Sequential Search."

[†]Free University Amsterdam, ICREA, IESE Business School, and University of Groningen, E-mail: jose.l.moraga@gmail.com.

[‡]University of Groningen, E-mail: z.sandor@rug.nl.

[§]Kelley School of Business, Indiana University, E-mail: mwildenb@indiana.edu.

1 Introduction

A significant body of work in economics has shown that search costs have far-reaching effects in economic activity. Well-known facts are that search costs alone can lead to price dispersion (Burdett and Judd, 1983; Stahl, 1989; Varian, 1980) as well as to wage and technology dispersion (Burdett and Mortensen, 1998; Acemoglu and Shimer, 2000). Search costs can also generate excessive product diversity in differentiated product markets (Wolinsky, 1984; Anderson and Renault, 2000) as well as inefficient quality investments (Wolinsky, 2005).

As a result, the estimation of consumer search costs has become an important area of empirical research. Hong and Shum (2006) were the first to develop a structural method to retrieve information on consumer search costs using market data. They focus on markets for homogeneous goods and present various approaches to estimate non-sequential and sequential consumer search models using only price data. Moraga-González and Wildenbeest (2008) present an alternative estimator based on maximum likelihood for non-sequential consumer search models. Hortaçsu and Syverson (2004) and Wildenbeest (2011) study search models where search frictions coexist with vertical product differentiation. In all these models, consumer search costs are found to be sizable and therefore an important source of market power.

The present paper studies the non-parametric identification and estimation of the costs of non-sequential search in markets for homogeneous products. It adds to the literature in three ways. First, we provide a proof that the critical search costs estimated in earlier work (cf. Hong and Shum, 2006; Moraga-González and Wildenbeest, 2008) are indeed non-parametrically identified. Second, we provide a new method based on semi-nonparametric (SNP) estimation that allows us to pool price data from different consumer markets with the same underlying search cost distribution but different valuations, selling costs, and numbers of competitors. Pooling data from different markets increases the number of estimated critical search cost cutoffs at all quantiles of the search cost distribution, which increases the precision of the estimate of the search cost CDF. Third, we provide sufficient conditions under which this type of data allows for the distribution of search costs to be identified on its full support.

The new method outperforms the spline approximation methods employed earlier in the literature (Hortaçsu and Syverson, 2004; Hong and Shum, 2006; Moraga-González and Wilden-

beest, 2008) with this type of data. Instead of estimating the parameters of the price distribution market by market, which ignores the link between the different data sets, our semi-nonparametric approach takes the search cost density as a parameter of the likelihood function and exploits all the data at once when estimating the model. SNP density estimators use a flexible polynomial-type parametric function that can approximate arbitrarily closely a large class of sufficiently smooth density functions (Gallant and Nychka, 1987), which means we obtain an essentially nonparametric estimator of the search cost distribution common to all markets. A Monte Carlo study illustrates that our estimator performs well in relatively small samples and outperforms existing market-by-market estimation methods.

To illustrate how our method works with real-world data we apply the SNP estimation procedure to a dataset of online prices for ten notebook memory chips. Our estimates indicate that median search costs are close to \$5. Search costs are quite dispersed; the majority of consumers visits at most three online stores before buying, while only a small fraction of consumers searches more than four times. Similar findings have been reported in several other empirical studies (Moraga-González and Wildenbeest, 2008; De los Santos, 2008; De los Santos *et al.*, 2011). Consumers with high search costs do not search much and therefore do not compare many prices, which gives substantial market power to the firms; as a result, estimated price-cost margins are significantly larger than what one would expect on the basis of the observed large number of firms operating in each market.

The structure of the paper is as follows. The next section reviews the non-sequential consumer search model. The identification result, the SNP estimation method, and the Monte Carlo study are presented in Section 3. Section 4 is devoted to the empirical application. Finally, Section 5 concludes. Our proofs are placed in the appendix to ease the reading.

2 The model

The model we study was proposed by Hong and Shum (2006) and generalizes the non-sequential consumer search model of Burdett and Judd (1983) by adding consumer search cost hetero-

generosity.¹ There is a finite number of firms K producing a good at constant returns to scale. Their identical unit cost is equal to r . There is a unit mass of buyers. Each consumer wishes to purchase a single unit of the good at most. The maximum price any buyer is willing to pay for the good is v . Consumers must engage in costly search to observe prices. We assume they search non-sequentially. In addition we assume that the first price quotation is observed at no cost.² Once a consumer has observed the desired number of prices, she chooses to buy from the store charging the lowest price. Consumers differ in their costs of search. A buyer's search cost is drawn independently from a common atomless distribution G with support $(0, \infty)$ and positive density g everywhere. A consumer with search cost c who searches k firms incurs a total search cost of $(k - 1)c$. The maximum number of prices a consumer can observe is K .

Firms and buyers play a simultaneous moves game. The solution concept is Nash equilibrium. An individual firm chooses its price strategy taking the price strategies of the rivals as well as consumers' search behavior as given. To allow for both pure and mixed pricing strategies, a firm i 's strategy is denoted by a probability distribution of prices F_i . Let F_{-i} denote the vector of pricing strategies used by firms other than i . The (expected) profit to firm i from charging a price p_i given the rivals' pricing strategies F_{-i} is denoted as $\Pi_i(p_i, F_{-i})$.

Likewise, an individual buyer takes as given the firms' pricing strategies and decides on her optimal search strategy to maximize her expected utility. The strategy of a consumer with search cost c is then a number k of prices to search for. Let the fraction of consumers searching k firms be denoted by μ_k .

We shall concentrate on symmetric Nash equilibria, i.e., equilibria where $F_i = F$ for all i . A symmetric equilibrium is a distribution of prices F and a collection $\{\mu_k\}_{k=1}^K$ such that:

- (a) $\Pi_i(p; F) \leq \bar{\Pi}$ for all p outside the support of F for all i ;
- (b) $\Pi_i(p; F)$ is equal to a constant $\bar{\Pi}$ for all p in the support of F , for all i ;
- (c) a consumer searching for the prices of k firms obtains no lower utility than by searching

¹Janssen and Moraga-González (2004) studied the same model with a search cost distribution with a two-point support.

²In our setting with search cost heterogeneity this assumption is inconsequential and can easily be relaxed at the cost of some additional notation. Earlier literature has assumed the first search to be costless in order to avoid problems of existence of equilibrium (cf. Diamond's (1971) paradox). To keep with the earlier literature we also maintain it here.

for any other number of prices;

$$(d) \sum_{k=1}^K \mu_k = 1.$$

Condition (a) is the standard Nash requirement that a firm must play a best-response to the strategies of the other players; condition (b) says that if the firms use a mixed strategy in equilibrium, then they must be indifferent among all the prices in the support of F ; finally, conditions (c) and (d) require the consumers to search such that their (expected) utility is maximized. Let us denote the equilibrium density of prices by f , with maximum price \bar{p} and minimum price \underline{p} .

2.1 Nash equilibria

We start by observing that in our game there may be two types of equilibria. There may be an equilibrium in pure strategies in which all firms charge a price equal to v and consumers optimally respond by not searching at all and visiting just one firm. This equilibrium is rejected right away in most empirical settings since we typically observe firms charging different prices while consumers are actively searching the market. Moreover, this equilibrium is non-robust in the sense that it heavily relies on the assumption that the first search is conducted at no cost; in fact, when the first search is costly this pure-strategy equilibrium fails to exist.

There may also be an equilibrium in mixed pricing strategies. In this equilibrium firm prices are dispersed and consumers respond by searching optimally to maximize their expected utility. Since both are common to most empirical settings this will be the equilibrium we will focus on. For this equilibrium to exist, there must be some consumers who search for one price only and others who search for two prices or more, i.e., $1 > \mu_1 > 0$ and $\mu_k > 0$ for some $k = 2, 3, \dots, K$. The intuition behind this observation is as follows. Suppose all consumers did search at least for two prices. If this were so, all firms would then be subject to price comparisons with rival firms. As a result, the firms would encounter themselves in a situation identical to the so-called Bertrand paradox (see e.g. Tirole, 1988). In such a situation, firms would have an incentive to undercut one another and thus all prices would be equal to marginal cost. Suppose now that no consumer did search at all. Then, since consumers would not be able to compare the prices

of different firms, they would charge the monopoly price v and again there would not be price dispersion in the market.

Our second observation is that, given that in equilibrium some consumers must search just once and others more than once, it must be the case that in symmetric equilibrium firms play mixed strategies with atomless price distributions.³ The intuition behind this remark is as follows. Suppose firms used distributions with an atom at a price $p \in (r, v]$. Since a price-tie at p would occur with strictly positive probability, an individual firm would gain by undercutting the tied price p , thereby attracting a larger share of the consumers who search for various prices and so obtaining greater profits. If there is an atom at $p = r$, then an individual firm would obtain zero profits; because some consumers do not search at all and therefore accept any price below or at v , the firm would have an incentive to deviate by increasing its price.

A third remark is that the upper bound of the symmetric equilibrium price distribution F must be equal to v . The reason for this is as follows. Suppose the upper bound were $\bar{p} < v$ and consider a firm charging \bar{p} . Since this firm would not sell to any of the consumers who search for various prices, its payoff would simply be equal to $(\bar{p} - r)\mu_1/K$, which is strictly increasing in \bar{p} ; as a result the firm would gain by deviating and charging v instead of \bar{p} .

Now that we have presented the basic properties of the mixed pricing strategy of the firms, let us consider the problem faced by the consumers. A consumer with search cost c must choose a number of prices to maximize her expected utility, where expected utility is equal to the difference between the consumer's valuation and the price she expects to pay, minus the cost of searching. If the consumer picks k prices to be searched, her expected utility is therefore given by $v - Ep_{1:k} - (k - 1)c$, where $Ep_{1:k}$ is short-hand notation for $E[\min\{p_1, p_2, \dots, p_k\}]$ and E indicates the expectation operator. Since every price is a random draw from F , the distribution of the minimum of k prices is equal to $(1 - F(p))^k$. Therefore, the number of prices that maximizes the utility of a consumer with search cost c is given by

$$k(c) = \arg \min_k \left[(k - 1)c + \int_p^v pk(1 - F(p))^{k-1} f(p) dp \right]. \quad (1)$$

Since $k(c)$ must be an integer, the problem in equation (1) induces a partition of the set of

³That is, discrete distributions or continuous distributions with "jumps" can be ruled out.

consumers into groups μ_k of consumers searching for k prices, with the property that $\sum_{k=1}^K \mu_k = 1$. We now describe such a partition. First we define the search cost cutoff c_k as the search cost of a consumer indifferent between searching for k prices, which gives her a utility equal to $v - Ep_{1:k} - (k - 1)c$, and searching for $k + 1$ prices, which allows her to obtain a utility level equal to $v - Ep_{1:k} - (k - 1)c$. Solving for c_k gives

$$c_k = Ep_{1:k} - Ep_{1:k+1}, \quad k = 1, 2, \dots, K - 1, \quad (2)$$

Since c_k decreases in k ,⁴ the fractions of consumers searching for k prices, denoted by μ_k , are given by

$$\mu_1 = 1 - G(c_1); \quad (3a)$$

$$\mu_k = G(c_{k-1}) - G(c_k), \quad k = 2, 3, \dots, K. \quad (3b)$$

To complete the equilibrium characterization, we now look at how firms should choose the distribution of prices F to maximize their profits given consumers' search behavior. The expected profit to a firm i from charging price p_i when rivals are setting prices according to the pricing strategy F is given by

$$\Pi_i(p_i; F) = (p_i - r) \left[\sum_{k=1}^K \frac{k}{K} \mu_k (1 - F(p_i))^{k-1} \right],$$

To understand this equation, note that firm i obtains a per consumer profit of $p_i - r$ and sells to a consumer who searches for k prices whenever the prices of the other $k - 1$ firms observed by this consumer are all higher than the price of firm i , which occurs with probability $(1 - F(p_i))^{k-1}$.

In a mixed strategy equilibrium, all the prices in the support of F must give the firm the same level of profits. Thus, for any price p in the support of F it must be the case that

⁴The cutoffs $c_k = Ep_{1:k} - Ep_{1:k+1}$ are in fact strictly monotonically decreasing in k because $Ep_{1:k}$ is strictly convex in k . A proof of this is available from the authors upon request. See also Stigler (1961), who mentions this property.

$\Pi_i(p; F) = \Pi_i(v; F)$. As a result, equilibrium requires

$$(p - r) \left[\sum_{k=1}^K k \mu_k (1 - F(p))^{k-1} \right] = \mu_1 (v - r) \quad (4)$$

to hold for all prices p in the support of F . Setting $F = 0$ in this equation and solving for p gives the minimum price charged in the market:

$$\underline{p} = \frac{\mu_1 (v - r)}{\sum_{k=1}^K k \mu_k} + r. \quad (5)$$

In Moraga-González *et al.* (2010) we show that an equilibrium always exists.

3 Econometric analysis

The econometric problem is to estimate the search cost distribution G using price data. Hong and Shum (2006) and Moraga-González and Wildenbeest (2008) propose different methods that exploit equations (2) to (5) to estimate the search cost CDF. In what follows, we briefly explain the two methods proposed so far (for details we refer the reader to the original contributions of these authors).

Hong and Shum (2006) formulate the estimation of the unknown search cost distribution as a two-step procedure. They propose to estimate first the parameters $\{\mu_k\}_{k=1}^K$ of the equilibrium price distribution obtained from equation (4) by maximum empirical likelihood (MEL), and then to recover the collection of cutoffs in equation (2) using the empirical CDF of prices. Suppose the researcher has a (large) data set with n prices and suppose $K (\leq n - 1)$ is the maximum number of prices a consumer may observe in the market. Let us assume each price p_j has probability π_j . Using equilibrium condition (4), for each price p_i we have the approximate equality

$$(p_i - r) \left[\sum_{k=1}^K k \mu_k \left(1 - \left[\sum_{j=1}^n \pi_j \mathbf{1}(p_j \leq p_i) \right] \right)^{k-1} \right] \simeq (v - r) \mu_1, \quad (6)$$

which can be transformed into a number $Q \geq K$ of population quantile restrictions:

$$\sum_{j=1}^n \pi_j \left[\mathbf{1} \left(p_j \leq r + \frac{(v-r)\mu_1}{\sum_{k=1}^K k\mu_k (1-s_\ell)^{k-1}} \right) - s_\ell \right] \simeq 0 \quad (7)$$

for $s_\ell \in [0, 1]$, $\ell = 1, 2, \dots, Q$. Using the lower bound defined in equation (5) one can eliminate marginal cost r from these constraints. Then, using MEL based on these constraints, one can obtain estimates of the parameters $\{\mu_k\}_{k=1}^K$. Finally, by combining these estimates with the cutoff points in equation(2) obtained directly from the empirical CDF of prices, one gets K points $\{(c_k, G(c_k))\}_{k=1}^K$ on the search cost distribution. These points serve to construct an estimate of the search cost CDF by interpolation.

Moraga-González and Wildenbeest (2008) put forward an alternative maximum likelihood (ML) method. There are two differences with respect to Hong and Shum's method. First, they compute the likelihood of a price as a function of the distribution of prices and exploit the equilibrium constancy-of-profits condition (4) to numerically calculate the value of the price CDF. In this way they obtain ML estimates of the parameters $\{\mu_k\}_{k=1}^K$. The second difference is that they introduce a method to compute ML estimates of the cutoffs by rewriting equation (2) as

$$c_k = \int_0^1 p(z)[(k+1)z - 1](1-z)^{k-1} dz, \quad k = 1, 2, \dots, K-1. \quad (8)$$

where $p(z)$ is the inverse of the price distribution obtained from equation (4):

$$p(z) = \frac{\mu_1(v-r)}{\sum_{k=1}^K k\mu_k(1-z)^{k-1}} + r. \quad (9)$$

These two methods yield estimates of the points $\{(c_k, G(c_k))\}_{k=1}^K$ of the search cost distribution. Under the standard regularity conditions, these points are estimated consistently. These two papers base their asymptotics on the number of prices n going to infinity. Although one of the regularity conditions is identification of the points $\{(c_k, G(c_k))\}_{k=1}^K$ of the search cost CDF, none of the earlier papers studied the identification issue. In the next subsection, we show that the sequence of points $\{(c_k, G(c_k))\}_{k=1}^K$ is identified.

3.1 Identification of search costs

In this subsection we ask whether the search cost distribution can be non-parametrically identified when the price distribution is known by the researcher. This treatment of the identification problem is common in nonparametric estimation and is in the spirit of Koopmans and Reiersøl (1950).

The analysis in Section 2 shows how the model maps the search cost distribution into the equilibrium distribution of prices. A feature of the model is that the entire set of consumers is partitioned into K groups of them. As a result, using the price equilibrium mapping, one can only hope to recover the (countable) sequence of points $\{(c_k, G(c_k))\}_{k=1}^K$ of the search cost CDF. The proposition below, proved in Appendix A, shows indeed that if we know the price distribution F , then we can identify the values of the search cost CDF corresponding to the cutoffs $\{c_k\}_{k=1}^K$.

Proposition 1 *Suppose that the econometrician observes the equilibrium price distribution F with support (\underline{p}, v) , which is continuous and is generated by the vector of variables (G, v, r, K) through equations (2), (3a), (3b), and (4). Then the points of the search cost distribution G corresponding to the sequence $\{c_k\}_{k=1}^K$ are identified.*

Obviously, when K is small, the sequence of points $\{(c_k, G(c_k))\}_{k=1}^K$ will be insufficient to obtain a precise estimate of the search cost distribution. The question that arises is whether such sequence allows us to recover the search cost CDF when $K \rightarrow \infty$. As illustrated in Figure 1, even if K is very large the search cost cutoffs do not give much information on the magnitude of search costs at quantiles other than zero. In this figure we plot the critical cutoff points c_k for different K 's (in particular, $K = 10, 15, 50$, and 100). In these plots we set $v = 500$ and $r = 50$, and assume consumer search costs follow a log-normal distribution with parameters $(a, b) = (0.5, 5)$.

[Figure 1 about here.]

We view this as a problem of identification of the search cost distribution in the relevant support. The problem is that data from a single market does not provide the econometrician with sufficient information to recover search costs at relatively high quantiles. The purpose

of the remainder of this subsection is to deal with this problem. Our proposal consists of bringing additional information to be able to estimate search costs more fully. Pooling data from various markets with similar search technology but different valuations, selling costs or numbers of competitors naturally lends itself as a feasible strategy to solve this identification problem. Implementing this idea in practice is not straightforward and in Section 3.2 we propose a new estimation method to do it.

The next proposition presents sufficient conditions under which the search cost distribution is identified on its full support using price data from multiple markets.

Proposition 2 *Suppose that there are infinitely (countably) many markets, indexed by m , all of them with the same underlying search cost distribution G . In every market m , the econometrician observes the price distribution function F^m with support (\underline{p}^m, v^m) , which is continuous and is generated by the vector of variables (G, v^m, r^m, K^m) through equations (2), (3a), (3b), and (4). In addition, assume that the difference between valuations and marginal costs $\{v^m - r^m\}_{m \geq 1}$ are random variables drawn independently from a distribution with support $(0, \infty)$. Then G is identified on the interval $[0, \sup_m c_1^m]$, where $\sup_m c_1^m = \sup\{c_1^m : m = 1, 2, \dots\}$ is the supremum of the set of c_1 -cutoff points from all the markets.*

This result says that one can solve the identification problem mentioned above by combining price data from various markets for which the search technology is similar but there is variation in valuations, selling costs, and numbers of competitors. Gathering the appropriate data is then relatively easy for the researcher. One just needs to take markets for different products in which consumers search for low prices in a similar fashion. To provide an example, if one aims to estimate the costs of search in the market for carpentry, one could pool data from the various professional services needed to refurbish a house: a carpenter, an electrician, a painter, a plumber, a bricklayer, a tiler, etc. The search technology to find acceptable prices for these professional services is basically the same; however, valuations, costs and the number of available professionals can be quite different across these services. This sort of data will do. If alternatively one considers the costs of search for prices on the Internet, one could take markets for different books, CDs, or DVD movies; in our application in Section 4, we use data from multiple markets for memory chips.

Intuitively, pooling data from various markets solves the problem of identification because every market generates a distinctive sequence of cutoff points, $\{c_k^m\}_{k=1}^{K^m}$, and this forces the search cost distribution function to be uniquely determined for a larger set of points, $\{\{c_k^m\}_{k=1}^{K^m}\}_{m \geq 1}$. Under the (large support) condition of the proposition, this set of points can be shown to be dense in the interval $[0, \sup_m c_1^m]$. If $\sup_m c_1^m = \infty$, then this proposition establishes identification of the search cost distribution in the full support.

We note that the (large support) condition in Proposition 2 is a sufficient condition also used in related nonparametric identification problems (see e.g. Matzkin, 1992; Matzkin, 1993; Ichimura and Thompson, 1998; Berry and Haile, 2009). In our case, we have adopted it to simplify the proof of identification. The assumption allows us to rely only on the cutoff points c_1^m to show identification.⁵

[Figure 2 about here.]

The ideas put forward in Proposition 2 are illustrated in Figure 2, where we plot the critical cutoff points c_k obtained from using data from $M = 1, 5, 25,$ and 50 markets. For each of these markets, we assume the number of firms is 10. In these plots we set $r = 50$ and again assume consumer search costs follow a log-normal distribution with parameters $(a, b) = (0.5, 5)$. For the case of data from one market only we set $v^m = 500$. For the situation with M markets we take valuations in market m as follows: $v^m = 100 + (500 - 100)(m - 1)/M$, so the lowest consumer valuation is always 100 and if there are for example five markets we get $\{v^m\}_{m=1}^5 = \{100, 200, 300, 400, 500\}$. The graphs make it clear that by using data from multiple markets we obtain much more information on the magnitude of search costs at high quantiles.

3.2 Estimation of search costs

As mentioned above, the previous studies on estimation of search cost distributions proceed in three steps: first, the parameters $\{\mu_k\}_{k=1}^K$ of the price distribution are estimated; second, the

⁵Since the large support condition is a sufficient condition, the identification result could be true under a weaker assumption. The proof would however be much more difficult because one would need to use the additional variation obtained from the other cutoff points c_k^m 's. The problem is that the mathematical relationship among all the c_k 's in a market, which is given by system of equations (8), is highly nonlinear and therefore using the other cutoffs in the proof of Proposition 2 is quite difficult. We would also like to note that the variation in K^m across markets is another source of variation that is useful in applications.

search cost cutoff points $\{c_k\}_{k=1}^K$ are obtained using the parameters of the price distribution; finally, a spline approximation of the search cost distribution is constructed by interpolating the sequence of estimated points $\{(c_k, G(c_k))\}_{k=1}^K$. As shown in the previous section, identification requires to pool data from many markets and in such a framework this earlier procedure presents some problems. In fact, it has to be applied market by market, in which case the researcher obtains multiple search cost estimates, one for each market. Interpolation is no longer feasible and one would have to fit a curve through the M estimated search cost distributions. It is not clear how one should proceed in such a case. For example, because the number of competitors K^m varies from market to market, whether all the markets should be allocated the same weight when fitting the curve is unclear. These difficulties lead us to propose a new estimation method that addresses these issues naturally.

We propose to employ semi-nonparametric (SNP) maximum likelihood estimation (Gallant and Nychka, 1987). The advantage of this method is that it is not applied market by market but designed to maximize the likelihood from all the markets jointly. In this way, the SNP procedure exploits the link between the prices not only within a market but also across markets because they all have the same underlying search cost distribution. We note that this method is different in essence because it takes the search cost density directly as the parameter of the likelihood. In this sense, it exploits the data more efficiently than the previous spline methods, since those rely on estimating the parameters of the price distribution in every market separately and, therefore, ignore the link between the different data sets.

The idea behind SNP estimation is to use a flexible functional approximation of the search cost density. This functional approximation depends on a finite set of parameters to be estimated and this set can be made arbitrarily large as the number of observations goes to infinity. We construct our estimator of the search cost density by employing a flexible polynomial-type approximation, following the SNP estimation technique developed by Gallant and Nychka (1987).

The likelihood function can be constructed by deriving the density of prices in each market $m = 1, 2, \dots, M$ as a function of the search cost density g . Let $f^m(p_i|g)$ denote the density of price p_i observed in a market m given the search cost distribution g . Since the prices

in a market m are independent draws from F^m , the log-likelihood function in market m is $L^m(g|\mathbf{p}_m) = \sum_{i=1}^{K^m} \log f^m(p_i|g)$ where \mathbf{p}_m is the K^m -dimensional vector of prices in market m . In order to compute this, first we apply the implicit function theorem to equation (4), which yields:

$$f^m(p_i|g) = \frac{\sum_{k=1}^{K^m} k \mu_k^m (1 - F^m(p_i|g))^{k-1}}{(p_i - r^m) \sum_{k=1}^{K^m} k(k-1) \mu_k^m (1 - F^m(p_i|g))^{k-2}}. \quad (10)$$

The quantities μ_k^m and r^m in this expression need to be computed in terms of g . By solving equation (5) for r^m we obtain an expression for the marginal cost in market m

$$r^m = \frac{\underline{p}^m \sum_{k=1}^{K^m} k \mu_k^m - \mu_1^m v^m}{\sum_{k=2}^{K^m} k \mu_k^m}. \quad (11)$$

We can (superconsistently) estimate the lower and upper bounds \underline{p}^m and v^m of the price distribution in a market m by taking the minimum and maximum prices, respectively, observed in the data.⁶ Then, for every market m , we compute $\{c_k^m\}_{k=1}^{K^m}$ from equations (8) in terms of g ,⁷ and then use equations (3a), (3b), (4), and (10) to find the values of $F^m(p_i|g)$ and $f^m(p_i|g)$. In this way we obtain the joint log-likelihood of all markets as a function of g :

$$L_M(g|\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M) = \frac{1}{M} \sum_{m=1}^M \left(\frac{1}{K^m} \sum_{i=1}^{K^m} \log f(p_i^m|g) \right)$$

For the polynomial-type parametric function that estimates the search cost density we employ the SNP density estimator of Gallant and Nychka (1987). This SNP estimator is based upon a Hermite polynomial expansion. The idea behind their SNP procedure is that any reasonable density can be mimicked by such a Hermite polynomial series. SNP density estimators are essentially nonparametric because the set of all Hermite polynomial expansions is dense in the set of density functions that are relevant (Gallant and Nychka, 1987).⁸

⁶In a similar fashion, order statistics are also used to estimate the lower and upper bound of distributions of bids (see e.g. Donald and Paarsch, 1993).

⁷In markets with many cutoff points solving this nonlinear system of equations may be time consuming. One alternative (used in the application) is to estimate the cutoffs directly by the empirical price CDF. The trade-off is precision of the estimates against computational time.

⁸SNP has recently been applied to the estimation of labor search frictions (Koning *et al.*, 2000), labor supply (Van Soest *et al.*, 2002), travel demand (Van der Klauw and Koning, 2003), and auctions (Brendstrup and Paarsch, 2006).

To apply the SNP estimation to our problem, we specify the search cost density as follows:

$$g(c; \gamma, \sigma, \theta) = \frac{\left[\sum_{i=0}^N \theta_i u_i(c) \right]^2}{\sum_{i=0}^N \theta_i^2}, \quad \theta \in \Theta_N \quad (12)$$

where $\Theta_N = \{\theta : \theta = (\theta_0, \theta_1, \dots, \theta_N), \theta_0 = 1\}$, N is the number of polynomial terms and

$$\begin{aligned} u_0(c) &= (c\sigma\sqrt{2\pi})^{-1/2} e^{-((\log c - \gamma)/\sigma)^2/4}, \\ u_1(c) &= (c\sigma\sqrt{2\pi})^{-1/2} ((\log c - \gamma)/\sigma) e^{-((\log c - \gamma)/\sigma)^2/4}, \\ u_i(c) &= \left[((\log c - \gamma)/\sigma) u_{i-1}(c) - \sqrt{i-1} u_{i-2}(c) \right] / \sqrt{i}, \text{ for } i \geq 2. \end{aligned}$$

This parametric form corresponds to the univariate SNP estimator studied extensively by Fenton and Gallant (1996). Our expressions are obtained by transforming their random variable x with the density defined in their Section 4.3 into $c = \exp(\gamma + \sigma x)$. This transformation is useful in our case since search costs are positive. The vector of parameters to be estimated by maximum likelihood is $\{\gamma, \sigma, \theta_1, \dots, \theta_N\}$ and N can be made arbitrarily large as the number of observations increases to infinity.

Gallant and Nychka (1987) provide conditions on the unknown density (e.g. differentiability and restricted tail behavior) under which their estimator is consistent using i.i.d. observations. In Appendix B we give details on how those conditions can be adapted to our search cost density.

In practice the number of polynomial terms N has to be chosen in an optimal way. For this, we can build on the cross-validation method of Coppejans and Gallant (2002). The essence of their cross-validation is to determine N for the data at hand by minimizing some loss function. Let f denote the true price density function and \hat{f}_N the price density function estimate computed as

$$\hat{f}_N(p) = f(p|\hat{g}_N),$$

where \hat{g}_N is the estimated search cost density with N polynomial terms. A standard way of

choosing N is by minimizing the integrated squared error (ISE), which in our case is

$$\int_{\underline{p}}^{\bar{p}} \left(\widehat{f}_N(p) - f(p) \right)^2 dp.$$

Since the true distribution $f(p)$ is not known, the ISE needs to be approximated.

There are various problem-specific ways to approximate the ISE; we proceed as follows.

First write

$$\int_{\underline{p}}^{\bar{p}} \left(\widehat{f}_N(p) - f(p) \right)^2 dp = \int_{\underline{p}}^{\bar{p}} \widehat{f}_N^2(p) dp - 2 \int_{\underline{p}}^{\bar{p}} \widehat{f}_N(p) f(p) dp + \int_{\underline{p}}^{\bar{p}} f^2(p) dp$$

and note that on the RHS only the first two terms depend on N . The first term $\int_{\underline{p}}^{\bar{p}} \widehat{f}_N^2(p) dp$ can be estimated (for example) by Monte Carlo simulations by drawing a sample from the uniform distribution on $[\underline{p}, \bar{p}]$. The integral from the second term can be written as

$$\int_{\underline{p}}^{\bar{p}} \widehat{f}_N(p) f(p) dp = E_P \left[\widehat{f}_N(p) \right],$$

which can be estimated by using the price observations in one market, i.e.,

$$\int_{\underline{p}}^{\bar{p}} \widehat{f}_N(p) f(p) dp \approx \frac{1}{K} \sum_{k=1}^K \widehat{f}_N(p_k).$$

In the empirical example the prices from different markets have different distributions, so we take the approximation of the average ISE across markets

$$\frac{1}{M} \sum_{m=1}^M \int_{\underline{p}}^{\bar{p}} \left(\widehat{f}_N^m(p) - f^m(p) \right)^2 dp.$$

To assess the performance of this method, we investigate it within the context of our Monte Carlo simulations in the next section.

3.3 Monte Carlo study

The purpose of the Monte Carlo study is threefold: (i) we study the small sample properties of the estimator; (ii) we study how using cross-validation to pick the number of polynomial terms

in the SNP density function performs in our setting; and (iii) we compare the performance of our estimator to an estimator that does not directly link different markets but instead estimates search costs market-by-market (based on Moraga-González and Wildenbeest, 2008). We focus on the estimation of the following search cost density:

$$g_0(c) = 0.5 \cdot \text{lognormal}(c, 2, 10) + 0.5 \cdot \text{lognormal}(c, 3, 0.5), \quad (13)$$

where $\text{lognormal}(c, a, b)$ refers to the densities of the lognormal distribution with parameters a and b , respectively. To make sure we are working in an environment that is not very different from the one used in our application in Section 4 we take $M = 10$ markets. Each market has the same search cost distribution but a different valuation net of marginal cost, $v^m - r^m$. The 10 values we use for $v^m - r^m$ are $\{40, 80, \dots, 400\}$. For each market m , we set K^m , the maximum number of prices a consumer can observe, equal to 35.⁹ With the parameters of a market m at hand, we compute the market equilibrium by numerically solving the system of equations (8). Given the cutoff values for a market m , we construct the equilibrium price distribution in that market m using equation (9). Next, we randomly draw 35 prices from each equilibrium price distribution F^m and use all 350 prices as an input for the SNP estimation procedure. The estimation is replicated 100 times.¹⁰

Table 1 shows the outcome of the Monte Carlo simulations for various values of N . The approximated (feasible) estimate of ISE selects $N = 8$, while the true (unfeasible) criterion selects $N = 6$ as the optimal number of polynomial terms. Search costs, however, are closest to the true search cost distribution when using $N = 8$, as shown in the last column of Table 1. This suggests our method works reasonably well.

[Table 1 about here.]

[Figure 3 about here.]

⁹Typically, the number of firms operating in a market will vary from market to market. Though this constitutes an additional source of variation, we do not need to use it here since we are assuming that the valuation net of marginal cost is different across markets.

¹⁰To gain computing time we use the empirical distribution of prices in each market to estimate the c_k 's. Although consistency of the estimator is preserved, this is likely to lead to less precise estimates, so our results should be seen as a lower bound on the performance of the estimator when using equation (8) instead.

[Figure 4 about here.]

Figures 3(a) and 3(b) show the estimated search cost distribution for $N = 8$. We report the mean and the 90 percent confidence interval of the 100 replications. Figure 3(a) corresponds to the search cost CDF, while Figure 3(b) corresponds to the search cost PDF. In both graphs, the solid curve represents the true search cost distribution, while the thick dashed curve shows the mean of the 100 estimations. The 90 percent confidence interval is given by the shaded area between the thin dashed curves. In spite of the relatively small number of markets and observations per market, the figures illustrate that our estimation procedure performs fairly well. The estimates mimic the true shape of the search cost CDF as well as PDF relatively well at most quantiles. Note that if we were to add more markets with relatively high valuation to our data set the number of search cost cutoffs would increase, which would improve the outcome of the estimation.

Existing approaches to estimate search costs (e.g., Hong and Shum, 2006; Moraga-González and Wildenbeest, 2008) are designed to estimate search costs market-by-market, while our SNP estimation procedure is specifically set up to maximize the joint likelihood from all markets. Figure 4(a) shows the estimated search cost PDF when we take the existing approach and use data for only one market.¹¹ Not only are the differences between the true search cost distribution (solid curves) and the mean of the 100 fitted distributions (thick dashed curves) larger than when using our multi-market SNP estimation procedure, also the 90 percent confidence interval (shaded area) is much wider. The search costs ISE confirms our visual findings: when taking data from just one market, the ISE takes on value 0.294×10^{-4} , which is almost four times as large as the corresponding ISE value for our SNP estimation procedure. If, alternatively, we use the data from all the markets and after estimating search costs market-by-market we take the average search cost density as an estimate of the overall search cost distribution, our SNP estimation procedure still outperforms the market-by-market approach, as illustrated in Figure 4(b). Although the search cost ISE in this case is slightly lowered to 0.229×10^{-4} , the 90 percent confidence interval widens. In sum, Figures 4(a) and 4(b) provide evidence that

¹¹We use prices for the market with $v = 400$ to make sure the maximum identifiable search cost value is the same as in our main specification. To obtain a parametric estimate of the search cost density we fit a SNP density function with $N = 8$ polynomial terms through the identified points on the search cost distribution, which are obtained using the approach in Moraga-González and Wildenbeest (2008).

the market-by-market approach underperforms vis-à-vis our multi-market SNP approach. It is less efficient because search costs are only constrained to be similar across markets after search costs have already been estimated for each market separately. However, we note that since the market-by-market approach is designed to maximize the likelihood function in each market separately it does an equally good job in terms of fitting the model to observed prices.

4 Application

In this section we use the SNP estimation method described above to quantify search costs in real-world markets for memory chips. We focus on computer memory chips for notebooks (so called SO-DIMM, or Small Outline Dual In-line Memory Module). Since we need products from different markets, we select memory chips produced for different brands and types of notebooks. Table 2 gives the details of the ten products we include in our data set. There are several reasons for choosing these memory chip data for the analysis. First, since all the chips are sold online, we expect search costs to be similar across markets. Second, even though all memory chips are manufactured by Kingston—the largest producer in the sector—each memory chip in our sample is meant to be used in a particular notebook brand only—including Toshiba, Dell, Acer, IBM and HP Compaq. Given that substitutability across products is somewhat limited due to technical reasons, we shall assume that different microchips belong in separate markets so the use of a search model with homogeneous products is reasonable.¹² All the memory chips we consider were somewhat at the top of the product line at the time of data collection. In particular, they exhibit relatively large storage capacity (1 gigabyte) and fast speed of operation (most of them above 400 MHz). Given the large storage capacity of the memory chips in the data set, most consumers would only consider to buy one memory chip, so the single-unit inelastic demand assumption of the theoretical model seems also reasonable.

[Table 2 about here.]

¹²Note that even though (within a market) the memory chips are exactly the same, the stores selling the chips might differ in terms of offered service, speed and quality of shipment, payment methods, etc. We come back to this issue at the end of this section.

For all the memory chips in the data set we collected online prices charged in the United States, in February 2006. To obtain a sufficiently representative sample, we gathered product and price information from several sources at the same time. We proceeded as follows. We first visited the price comparison sites *shopper.com* and *pricegrabber.com* and collected the names of all the shops that were seen active in markets for memory chips; in total we found 49 stores. If for a particular product we saw a shop quoting its price on *shopper.com* and/or *pricegrabber.com*, we took the price directly from the price comparison site; otherwise we visited the web-address of the vendor to check if the product was available and at what price it was offered.

Table 3 gives some summary statistics of the data set. The number of firms quoting prices in each market is relatively large, ranging from 24 to 41. In our study we estimate the maximum number of prices consumers can search for in each market K^m by the number of firms that were observed to be quoting prices in that market. Almost all memory chips are priced above \$100. For all products we observe significant price dispersion as measured by the price range (difference between the maximum and the minimum prices) and by the coefficient of variation. We note that the (gross) benefits to a consumer from searching are significant; in particular, the (gross) gains from searching all the firms thereby becoming fully informed relative to searching for one price only in these markets range from \$16.32 to \$33.05. As mentioned above, we estimate the valuation of a memory chip by the maximum price observed in the market.

[Table 3 about here.]

Our model assumes consumers search non-sequentially. Consumers obviously visit stores sequentially in the real world, so what truly distinguishes non-sequential search from sequential search is how consumers select the stores they visit—if they search non-sequentially, the number of stores searched is determined before searching, while if they search sequentially, the number of searches depends on what has been observed. Although non-sequential search is often thought of as a constrained version of sequential search, Morgan and Manning (1985) have shown that the optimal search rule is hybrid in nature: it includes decisions on the sample size as well as whether to continue searching or not. This means both non-sequential and sequential search are special cases—non-sequential search is typically optimal when the search outcome is observed with delay, for instance when applying for a job or college, or when obtaining estimates from

contractors. Even though a typical online shopper does not face such a delay when searching online, according to Manning and Morgan (1982) sufficiently large economies of scale when searching can make it optimal to search multiple firms at once without using the option to continue searching. This is typically the case when searching online: once having found the correct memory chip and price at one web store, simply copying and pasting the chip’s model number to another online store is all it takes to obtain an additional price quote. This might also explain why, using individual specific observations on browsing history, De los Santos *et al.* (2011) find that observed search patterns for online books are more consistent with non-sequential search than sequential search.

Because we only observe the stores’ prices at one moment in time, we cannot check whether stores indeed use mixed pricing strategies, as predicted by our search model. However, using a different data set Moraga-González and Wildenbeest (2008) show that firms indeed seem to mix prices in the online market for memory chips; at the same time, other studies find evidence for mixed strategies in other markets (e.g., Lach (2002) for chicken, refrigerators, coffee, and flour in Israel; Lewis (2008) for gasoline; and Wildenbeest (2011) for grocery products in the United Kingdom).

[Table 4 about here.]

We follow the procedure explained in Section 3.2 and use cross-validation for choosing the number of polynomial terms N . Table 4 gives the SNP estimation results for different values of N . These results are obtained using the empirical price CDF in each market to calculate the c_k ’s.¹³ As can be seen in the table, up to $N = 20$ there is a steady improvement in ISE, while for larger N the improvement is very small. We therefore use the estimated parameters for $N = 20$ to derive our estimate of the search cost distribution. The solid curve in Figure 5(a) denotes the estimated search cost CDF, while the shaded area indicates the 95 percent confidence interval.¹⁴ The graph also shows how the estimated search cost cutoffs (gray dots)

¹³In cases when there are sufficiently many observations, as is the case in our data set, we can use the empirical distribution of prices in each market directly to estimate the c_k ’s. The gain in computing time is huge and the results for our data are very similar. See also Figure 7(a).

¹⁴Since standard errors of the parameter estimates are only meaningful in the case where the presented model is the true parametric model, we have obtained the confidence interval using bootstrapping. For each replication we draw 10 markets with replacement; the 95 percent confidence interval is obtained using the 2.5th and 97.5th percentile of 100 replications.

cover the support of the search cost distribution.

[Figure 5 about here.]

Using the estimates of the parameters of the SNP specification we can compute the mean, the median, and the standard deviation of the search cost distribution. The median consumer has a search cost equal to \$5.05. On average a consumer has a search cost value equal to \$8.70 and the standard deviation is \$7.35. It is also interesting to investigate the distribution of search intensities in these markets. Since each market has specific parameters, even though search costs are assumed to be similar, it is unlikely that consumer search behavior will be the same across markets. Table 5 shows that it is indeed the case that search intensities are different across markets. For example, in the market for the KTT3311A memory chip, 26 percent of consumers searches for one price only while in the market for the KTH-ZD8000A memory chip the share of consumers who searches once is 34 percent. Similarly, in the market for the KTT3311A chip, 24 percent of consumers searches for two prices, while in the market for the KTD-INSP8200 memory chip the share of consumers who searches for two prices is 62 percent. However, the share of consumers searching at most three times is more or less similar across markets; approximately 91 percent of the consumers have search cost above \$3.70 and search for at most three prices. Table 5 also illustrates that the group of consumers searching for the prices of between 4 and 15 firms is with percentages between 0 and 4 relatively small. About 8 percent of consumers search for prices thoroughly; they search for the prices of more than 15 stores, which means they have search costs less than 43 dollar cents. Figures 5(a) and 5(b) show that the consumers can roughly be divided into three groups: buyers who do not search, buyers who search for at most three prices and buyers who search for many prices in the market.

Our findings are in line with several other empirical studies; Moraga-González and Wildenbeest (2008) report similar results using a different estimation method and data set, while Wildenbeest (2011) finds that very few consumers visit an intermediate number of stores when searching for grocery products, even if quality differentiation is taken into account. Moreover, our result that consumers search very little is supported by the consumer-specific web browsing data for online bookstores used in De los Santos (2008) and De los Santos *et al.* (2011).

[Table 5 about here.]

The fact that a significant proportion of consumers does not search for many prices confers substantial market power to the firms. Using the estimates of the SNP specification, we can retrieve the marginal cost r in each market, which is also reported in Table 5. Marginal costs range between 56 and 64 percent of the value of the product, while average price-cost margins range between 19 and 24 percent across markets.

[Figure 6 about here.]

To test whether the estimated model explains observed prices well, we calculate the Kolmogorov-Smirnov statistic (KS-test) in each individual market. The KS-test statistic is based on the maximum difference between the empirical price CDF and the estimated price CDF, which is the computed price equilibrium given the estimate of the search cost distribution. The null hypothesis for this test is that the distributions are similar, the alternative hypothesis is that the empirical and the estimated price CDF are different. Table 5 gives the KS-test results—for eight out of ten memory chips the KS value is below the 95 percent critical value of the KS-statistic of 1.36, which means that for these chips we cannot reject the null-hypothesis that the prices are drawn from the estimated price CDF.¹⁵ The goodness-of-fit is also shown in Figure 6, where we have plotted both the empirical and the estimated price CDF for two of the ten markets. Figure 6(a) shows the fit for the memory chip that gives the best fit; the empirical price CDF, as indicated by the dashed line, is close to the estimated price CDF, which is represented by the solid curve. Also plotted is the band that gives the maximum allowed difference between the estimated and empirical price CDF. Figure 6(b) shows that the empirical price CDF is just outside this band for the memory chip that gives the worst fit.

We have estimated the model using the empirical price CDF in each market to calculate the search cost cutoffs. The main reason for doing so is the gain in computing time: this avoids having to solve the system of equations in equation (8) in each function evaluation. Figure 7(a) shows that the search cost CDF when the c_k 's are estimated (dashed curve) is not very different from the one obtained when using the empirical CDF to get the cutoffs (solid curve).

¹⁵We have calculated KS in Table 5 as $\sqrt{K^m} \cdot \tau_{K^m}$, where K^m is the number of price observations for the specific memory chip and τ_{K^m} is the maximum absolute difference over all prices between the estimated price CDF and the empirical price CDF.

The prices used for our estimations include neither shipping costs nor sales taxes. The main reason for leaving these out is that shipping costs and sales taxes depend on the state in which the consumer resides, which makes it difficult to compare total prices. However, for robustness purposes, we also estimate the model neglecting sales taxes but including shipping costs for residents of New York. Figure 7(b) shows that the estimated search cost CDF (dashed) is very similar to the search cost CDF obtained when ignoring shipping costs (solid).

[Figure 7 about here.]

Although the memory chips themselves are completely homogeneous, the price differences across vendors for a given chip may be due to store differentiation. Consumers might prefer one shop over another on the basis of observable store characteristics like quality ratings, return policies, stock availability, order fulfillment, payment methods, etc. To see the impact of observable shop characteristics on prices, we regress prices on indicators that are readily available from the price comparison sites. More precisely, we estimate the following model:

$$PRICE_j = \beta_0 + \beta_1 \cdot RATING_j + \beta_2 \cdot DISCLOSE_j + \beta_3 \cdot STOCK_j + \beta_4 \cdot LOGO_j + \varepsilon_j, \quad (14)$$

where, for each product, $PRICE_j$ is the list price of store j , $RATING_j$ is an average of the ranking of store j on shopper.com and pricegrabber.com, $DISCLOSE_j$ is a dummy for whether shop j disclosed shipping cost on either shopper.com or pricegrabber.com, $STOCK_j$ is a dummy for whether shop j had the item in stock, and $LOGO_j$ is a dummy for whether shop j had its logo on either shopper.com or pricegrabber.com. We estimate this equation by OLS. The resulting R-squared values indicate that only between 3 and 27 percent of the total variation in prices can be attributed to observable differences in store characteristics.¹⁶ Although this does not rule out that there are unobservable differences between stores (e.g., cost differences or branding), this does suggest that the observable characteristics cannot explain the vast majority of variation in prices and that something else must cause such variability. In spite of this, for robustness purposes, we also estimate the model using the residuals of the regression above. This is standard practice in many structural auction models (e.g. Haile *et al.*, 2003; Bajari *et*

¹⁶For all memory chips, all the OLS coefficient estimates were not significant except the coefficient for $LOGO_j$, which was positive and significant at a 5 percent level for the KTM-TP3840 and KTH-ZD8000A chips.

al., 2006; An *et al.*, 2010); Wildenbeest (2011) shows that if stores obtain quality input factors in perfectly competitive markets, the quality production function exhibits constant returns to scale, and consumers have the same preferences towards quality, this procedure is theoretically correct (moreover, all our results on identification and consistency of the estimator hold for such specification as well). As shown in Figure 7(c), estimated search costs are uniformly lower. This result is intuitive: when taking store heterogeneity into account, the (gross) gains from searching will be lower, which means that in order to explain observed prices consumers should have lower search costs and search more than in the model without store heterogeneity.

Finally, Figure 7(d) shows the estimated search cost CDF when estimating search costs market-by-market using the approach put forward by Moraga-González and Wildenbeest (2008). This estimate is obtained by fitting an SNP density function to the estimated search costs in each market and taking the average. As can be seen from the graph, there are some differences with our main specification: while our SNP procedure predicts only 9 percent of consumers have search costs less than \$3.70, this would be 34 percent according to the market-by-market approach.

5 Conclusions

Since the seminal contribution of Stigler (1961), economists have dedicated a significant amount of effort to understand the nature of competition in markets where price information is not readily available to consumers. One of the lessons learnt is that consumer search models may lead to price dispersion, a prediction quite different from the ‘law of one price’ obtained from conventional economic theory. Another is that the particular direction of the effects of public policy measures such as the introduction of taxes or the dismantling of barriers to entry depends on the shape of the search cost distribution. These observations motivate the development of methods to estimate search costs to be used in the simulation of counterfactual scenarios. The estimation of consumer search costs is nowadays an important area of empirical research.

This paper has studied the non-parametric identification and estimation of the costs of simultaneous search in markets for homogeneous products. We have argued that in order to increase the precision of the estimate of the search cost distribution one needs to increase the

number of estimated critical search cost cutoffs in all quantiles of the search cost CDF. We have shown this can be done by pooling price data from various markets with similar search technology but different valuations, firms' costs and numbers of competitors. To take advantage of the relationship between the distinct markets we have proposed a new method to estimate the search cost density function by a semi-nonparametric density estimator whose parameters maximize the joint likelihood corresponding to all the markets. The paper has also illustrated the potential of our method by applying it to a data set of online prices for ten notebook memory chips. The estimates obtained suggest that the search cost density is essentially bimodal such that a large fraction of consumers searches for very few prices and a small fraction of consumers searches for a relatively large number of prices.

Along the way we have made several simplifying assumptions. One of the assumptions has been that, within a market, consumers have the same valuation. In future work, we would like to relax this assumption and study a framework where there is heterogeneity both in consumer valuations and search costs. One of the advantages of developing such a framework is that it would enable the econometrician to estimate the correlation between consumer valuations and search costs. Another simplifying assumption has been that firms have complete information about the costs of one another. Our model could be extended to a setting with private information about the marginal costs of production. Estimation of such a model would enable us to distinguish price dispersion due to marginal cost heterogeneity from price dispersion due to search costs. Finally, an important restriction of our model has been that we treat the different markets as completely separated. In more general settings, one would like to develop a model of product differentiation with search costs. We strongly believe the ideas developed in this paper can be applied to such markets. The first step is to develop a tractable model that incorporates strategic price dispersion together with product heterogeneity. This is work we will pursue in the future.

References

Acemogly D, Shimer R. 2000. Wage and technology dispersion. *Review of Economic Studies* **67**: 585–607.

- An Y, Hu Y, Shum M. 2010. Estimating first-price auction models with unknown number of bidders: a misclassification approach. *Journal of Econometrics* **157**: 328–341.
- Anderson SP, Renault R. 2000. Consumer information and firm pricing: negative externalities from improved information. *International Economic Review* **41**: 721–742.
- Bajari P, Houghton S, Tadelis, S. 2006. *Bidding for incomplete contracts: an empirical analysis*. NBER Working Paper 12051.
- Berry ST, Haile PA. 2009. *Nonparametric identification of multinomial choice demand models with heterogeneous consumers*. NBER Working Paper 15276.
- Brendstrup B, Paarsch HJ. 2006. Identification and estimation in sequential, asymmetric, English auctions. *Journal of Econometrics* **134**: 69–94.
- Burdett K, Judd KL. 1983. Equilibrium price dispersion. *Econometrica* **51**: 955–969.
- Burdett K, Mortensen DT. 1998. Wage differentials, employer size, and unemployment. *International Economic Review* **39**: 257–273.
- Coppejans M, Gallant AR. 2002. Cross-validated SNP density estimates. *Journal of Econometrics* **110**: 27–65.
- De los Santos B. 2008. *Consumer search on the Internet*. NET Institute Working Paper 08-15.
- De los Santos B, Hortacısu A, Wildenbeest MR. 2011. Testing models of consumer search using data on web browsing and purchasing behavior. *American Economic Review*, forthcoming.
- Diamond PA. 1971. A model of price adjustment. *Journal of Economic Theory* **3**: 156–168.
- Donald SG, Paarsch HJ. 1993. Piecewise pseudo-maximum likelihood estimation in empirical models of auctions. *International Economic Review* **34**: 121–148.
- Fenton VM, Gallant AR. 1996. Qualitative and asymptotic performance of SNP density estimators. *Journal of Econometrics* **74**: 77–118.
- Gallant AR, Nychka DW. 1987. Semi-nonparametric maximum likelihood estimation. *Econometrica* **55**: 363–390.

- Haile PA, Hong H, Shum M. 2003. *Nonparametric tests for common values in first-price sealed-bid auctions*. NBER Working Paper 10105.
- Hong H, Shum M. 2006. Using price distributions to estimate search costs. *RAND Journal of Economics* **37**: 257–275.
- Hortaçsu A, Syverson C. 2004. Product differentiation, search costs, and competition in the mutual fund industry: a case study of S&P 500 index funds. *Quarterly Journal of Economics* **119**: 403–456.
- Ichimura H, Thompson TS. 1998. Maximum likelihood estimation of a binary choice model with random coefficients of unknown distribution. *Journal of Econometrics* **86**: 269–295.
- Janssen MCW, Moraga-González JL. 2004. Strategic pricing, consumer search and the number of firms. *Review of Economic Studies* **71**: 1089–1118.
- van der Klaauw B, Koning RH. 2003. Testing the normality assumption in the sample selection model with an application to travel demand. *Journal of Business & Economic Statistics* **21**: 31–42.
- Koning P, van den Berg GJ, Ridder G. 2000. Semi-nonparametric estimation of an equilibrium search model. *Oxford Bulletin of Economics and Statistics* **62**: 327–356.
- Koopmans TC, Reiersøl O. 1950. The identification of structural characteristics. *Annals of Mathematical Statistics* **21**: 165–181.
- Lach S. 2002. Existence and persistence of price dispersion: an empirical analysis. *Review of Economics and Statistics* **84**: 433–444.
- Lewis MS. 2008. Price dispersion and competition with differentiated sellers. *Journal of Industrial Economics* **56**: 654–678.
- Manning R, Morgan P. 1982. Search and consumer theory. *Review of Economic Studies* **49**: 203–216.
- Matzkin RL. 1992. Nonparametric and distribution-free estimation of the binary choice and threshold crossing models. *Econometrica* **60**: 239–270.

- Matzkin RL. 1993. Nonparametric identification and estimation of polychotomous choice models. *Journal of Econometrics* **58**: 137–168.
- Moraga-González JL, Wildenbeest MR. 2008. Maximum likelihood estimation of search costs. *European Economic Review* **52**: 820–48.
- Moraga-González JL, Sándor Z, Wildenbeest MR. 2010. *Nonsequential search equilibrium with search cost heterogeneity*. IESE Business School Working Paper 869.
- Morgan P, Manning R. 1985. Optimal Search. *Econometrica* **53**: 923–944.
- van Soest A, Das M, Gong X. 2002. A Structural Labour Supply Model with Flexible Preferences. *Journal of Econometrics* **107**: 345–374.
- Stahl DO. 1989. Oligopolistic Pricing with Sequential Consumer Search. *American Economic Review* **79**: 700–712.
- Stigler G. 1961. The Economics of Information. *Journal of Political Economy* **69**: 213–225.
- Tirole, J. 1988. *The Theory of Industrial Organization*. The MIT Press.
- Varian HR. 1980. A Model of Sales. *American Economic Review* **70**: 651–659.
- Wald A. 1949. Note on the Consistency of the Maximum Likelihood Estimate. *Annals of Mathematical Statistics* **20**: 595–601.
- Wildenbeest MR. 2011. An Empirical Model of Search with Vertically Differentiated Products. *RAND Journal of Economics*, forthcoming.
- Wolinsky A. 1984. Product Differentiation with Imperfect Information. *Review of Economic Studies* **51**: 53–61.
- Wolinsky A. 2005. Procurement via Sequential Search. *Journal of Political Economy* **113**: 785–810.

A APPENDIX A

Proof of Proposition 1. Consider the triplets of variables $(F, \{\mu_k\}_{k=1}^K, \{c_k\}_{k=1}^K)$ and $(F, \{\mu'_k\}_{k=1}^K, \{c'_k\}_{k=1}^K)$ that are generated by the quadruplets of variables (G, v, r, K) and (G', v', r', K') , respectively, where G' is another distribution function with support $(0, \infty)$ and positive density on this support. Then we prove the result by showing that $\mu'_k = \mu_k$, $c'_k = c_k$ and $G'(c_k) = G(c_k)$ for any $k \in \{1, 2, \dots, K\}$.

First we note that neither μ_1 nor μ'_1 can be equal to zero. Indeed, if $\mu_1 = 0$ then by equation (4) $\sum_{k=2}^K k\mu_k (1 - F(p))^{k-1} = 0$ for any $p \in [p, v]$, which, due to the fact that F is increasing and continuous, can only happen if $\mu_k = 0$ for any $k \geq 2$. This is in contradiction with $\sum_{k=1}^K \mu_k = 1$, so $\mu_1 > 0$. Since exactly the same arguments apply to μ'_1 , we have shown that μ_1 and μ'_1 are strictly positive.

Next we prove that $r' = r$. Since $\mu_1 > 0$, equation (4) implies that F is strictly increasing on its support and hence invertible. By putting $p = F^{-1}(1 - z)$ in (4) for $\{\mu_k\}_{k=1}^K, r$ and $\{\mu'_k\}_{k=1}^K, r'$, we obtain that

$$\frac{\mu_1(v-r)}{\sum_{k=1}^K k\mu_k z^{k-1}} + r = \frac{\mu'_1(v-r')}{\sum_{k=1}^K k\mu'_k z^{k-1}} + r' \quad \text{for any } z \in [0, 1].$$

This implies that

$$\mu_1(v-r) \left(\sum_{k=1}^K k\mu'_k z^{k-1} \right) - \mu'_1(v-r') \left(\sum_{k=1}^K k\mu_k z^{k-1} \right) - (r' - r) \left(\sum_{k=1}^K k\mu_k z^{k-1} \right) \left(\sum_{k=1}^K k\mu'_k z^{k-1} \right) = 0$$

for any $z \in [0, 1]$. Since the LHS is a polynomial in z , all its coefficients must be equal to 0.

Suppose by contradiction that $r' \neq r$. This implies that

$$\text{either } \mu_K = \mu_{K-1} = \dots = \mu_2 = 0 \text{ (so } \mu_1 = 1) \text{ or } \mu'_K = \mu'_{K-1} = \dots = \mu'_2 = 0 \text{ (so } \mu'_1 = 1). \quad (\text{A15})$$

Indeed, by contradiction assume that (A15) does not hold; then let $M, M' \geq 2$ denote the maxima of k, ℓ such that $\mu_k > 0$ and $\mu'_\ell > 0$. The coefficient of $z^{M+M'-2}$ is $-(r' - r) M\mu_M M'\mu'_{M'}$, which must be equal to 0, a contradiction with our assumptions. Therefore, (A15) must hold. In either case we have a contradiction because (5) implies that $\underline{p} = v$, which means that the

price distribution F is degenerated. This establishes that $r' = r$.

Next we show that $\mu'_k = \mu_k$ for any k . From equation (4) we obtain

$$\sum_{k=1}^K k \frac{\mu_k}{\mu_1} (1 - F(p))^{k-1} = \frac{v-r}{p-r} = \sum_{k=1}^K k \frac{\mu'_k}{\mu'_1} (1 - F(p))^{k-1} \quad \text{for any } p \in [\underline{p}, v].$$

This is equivalent to

$$\sum_{k=2}^K k \left(\frac{\mu_k}{\mu_1} - \frac{\mu'_k}{\mu'_1} \right) z^{k-1} = 0 \quad \text{for any } z \in [0, 1]. \quad (\text{A16})$$

Since the LHS is a polynomial in z , all its coefficients must be equal to 0. Therefore, $\frac{\mu_k}{\mu_1} = \frac{\mu'_k}{\mu'_1}$ for $k = 2, \dots, K$. On the other hand, $\mu_1 + \sum_{k \geq 2} \mu_k = \mu'_1 + \sum_{k \geq 2} \mu'_k = 1$. These equalities together imply $\frac{1}{\mu_1} = \frac{1}{\mu'_1}$ and therefore $\mu'_k = \mu_k$ for any $k \geq 1$.

The equalities $c'_k = c_k$ follow from equation (2). It remains to show that $G'(c_k) = G(c_k)$ for any $k \geq 1$. We do so by showing that $\{G(c_k)\}_{k \geq 1}$ is uniquely determined by the series $\{\mu_k\}_{k \geq 1}$. By equations (3a) and (3b), $G(c_{k-1}) - G(c_k) = \mu_k$ for any $k \geq 1$. This implies that $G(c_k) = 1 - \sum_{h=1}^k \mu_h$ for any $k \geq 1$. The result then follows from the equality $\mu'_k = \mu_k$ for any $k \geq 1$ established above. ■

Proof of Proposition 2. In the proof we write $c_1(\theta)$ to make explicit the dependence of c_1 on $\theta \equiv v - r$. Note that due to the continuity of $c_1(\theta)$, $\sup_{\theta \in (0, \infty)} c_1(\theta) = \sup_m c_1^m$. Take an arbitrary interval $(a, b) \subset (0, \sup_{\theta \in (0, \infty)} c_1(\theta))$. Then the pre-image set defined as $c_1^{-1}(a, b) = \{\theta : c_1(\theta) \in (a, b)\}$ is a nonempty set, open in $(0, \infty)$ because $\lim_{\theta \rightarrow 0^+} c_1(\theta) = 0$ (by equation (4), if $\theta = 0$ then $\mu_k = 0$ for $k \geq 2$, so $\mu_1 = 1$ and thus $G(c_1) = 0$) and c_1 is, by assumption, a continuous function of θ . Therefore, with probability 1 there exists an m such that $\theta_m = v^m - r^m \in c_1^{-1}(a, b)$, which means that $c_1(\theta_m) \in (a, b)$.¹⁷ Because the interval (a, b) has been chosen arbitrarily, we have proven that for any interval, we can find an m such that the corresponding cutoff point $c_1(\theta_m)$ is included in the interval with probability 1. Since $G(\theta_m) = G(c_1(v^m - r^m))$, $m \geq 1$, are identified, this establishes that in an arbitrary interval

¹⁷The argument for this statement is the following. Suppose that we have iid random variables x_1, x_2, \dots, x_n drawn from a distribution with support $(0, \infty)$ and let $(c, d) \subset (0, \infty)$. Then the probability that at least one of these random variables is in (c, d) is equal to $1 - P(x_i \notin (c, d))^n = 1 - [1 - P(x_i \in (c, d))]^n$. Since $P(x_i \in (c, d)) > 0$, the above probability goes to 1 when $n \rightarrow \infty$. So when we have a countably infinite sequence of random variables, the probability that at least one of these random variables is in (c, d) is 1.

$(a, b) \subset (0, \sup_{\theta \in (0, \infty)} c_1(\theta))$ we can find a point at which the search cost distribution is identified with probability 1. Therefore, since it is continuous, G is identified on $[0, \sup_{\theta \in (0, \infty)} c_1(\theta)]$. ■

B APPENDIX B

In this section of the Appendix we adapt the general conditions in Gallant and Nychka (1987, henceforth GN) for the consistency of the search cost density estimator and discuss some primitive conditions specific to our model. Since the price observations in our model come from multiple markets that may be heterogeneous in valuations, firms' costs and number of firms, the price observations may not be i.i.d. In order to be able to treat the prices as i.i.d., we will regard these conditioning variables as random. This is not restrictive since it is just a matter of interpretation; in fact it is analogous to treating the covariates in a regression as random, in order to have i.i.d. dependent variables.

For this purpose, let us first modify the notation of the price density to $f(p|g; v^m, r^m, K^m)$ in order to make explicit the dependence on valuations, firms' costs and number of firms. Then

$$L_M(g) = \frac{1}{M} \sum_{m=1}^M \left(\frac{1}{K^m} \sum_{i=1}^{K^m} \log f(p_i^m | g; v^m, r^m, K^m) \right),$$

is the log-likelihood presented above, where for simpler notation we ignore the price vectors on the LHS. We regard the triplets $(v^m, r^m, K^m)_{m=1}^M$ as an i.i.d. sample of random variables. Then by Kolmogorov's strong law of large numbers, $L_M(g) \xrightarrow[M \rightarrow \infty]{a.s.} L(g) \equiv E[\log f(p|g; v, r, K)]$, provided that $E[\log f(p|g; v, r, K)] < \infty$ (this condition will follow from Lemma A.1 below).¹⁸

In order to state sufficient conditions for the consistency of our search cost density estimator,

¹⁸Note that

$$L(g) = E \left[\frac{1}{M} \sum_{m=1}^M \left(\frac{1}{K^m} \sum_{i=1}^{K^m} \log f(p_i | g; v^m, r^m, K^m) \right) \right].$$

Indeed, since $\left(\frac{1}{K^m} \sum_{i=1}^{K^m} \log f(p_i | g; v^m, r^m, K^m) \right)_{m=1}^M$ is an i.i.d. sample, by the law of iterative expectation,

$$\begin{aligned} E \left[\frac{1}{M} \sum_{m=1}^M \left(\frac{1}{K^m} \sum_{i=1}^{K^m} \log f(p_i | g; v^m, r^m, K^m) \right) \right] &= E \left[\frac{1}{K^m} \sum_{i=1}^{K^m} E[\log f(p_i | g; v^m, r^m, K^m) | K^m] \right] \\ &= E[E[\log f(p_i | g; v^m, r^m, K^m) | K^m]], \end{aligned}$$

where the last equality holds because in each market m the prices $(p_i)_{i=1}^{K^m}$ are i.i.d.. Then by the law of iterative expectation the statement follows.

we introduce some further notation. Recall that the search costs c we consider are exponential transformations of the random variables x from GN, that is, $c = \exp(\gamma + \sigma x)$. The density of c is $g(c) = \frac{1}{\sigma c} h\left(\frac{\log c - \gamma}{\sigma}\right)$, where h denotes the density of x . Let

$$\mathcal{G} = \left\{ g : g(c) = \frac{1}{\sigma c} h\left(\frac{\log c - \gamma}{\sigma}\right), \gamma \in \mathbb{R}, \sigma > 0, h \in \mathcal{H} \right\}$$

denote the set of admissible search cost densities, where \mathcal{H} is the set of admissible densities defined by GN (p.369). For each $\gamma \in \mathbb{R}$, $\sigma > 0$ define the operator $\|\cdot\| : \mathcal{G} \rightarrow \mathbb{R}$ such that $\|g\| = \|h\|_{GN}$, where $\|\cdot\|_{GN}$ is the consistency norm from GN (p.371), and define the operator $T : \mathcal{H} \rightarrow \mathcal{G}$ with $T(h)(c) = \frac{1}{\sigma c} h\left(\frac{\log c - \gamma}{\sigma}\right)$. Then $\|\cdot\|$ is a norm on \mathcal{G} and T is a homeomorphism between the normed spaces $(\mathcal{H}, \|\cdot\|_{GN})$ and $(\mathcal{G}, \|\cdot\|)$.

Let $g_0 \in \mathcal{G}$ be the true search cost density and $\mathcal{G}_N = \{g_N(\cdot; \gamma, \sigma, \theta) : \gamma \in \mathbb{R}, \sigma > 0, \theta \in \Theta_N\}$ the space of SNP estimators, where $g_N(\cdot; \gamma, \sigma, \theta)$ is defined in (12). Denote the SNP estimator of g_0 by \hat{g} , let the number of observations be n .

Proposition A.1 *Under the following conditions:*

- (a) *Compactness: The closure of \mathcal{G} is compact,*
- (b) *Denseness: $\cup_{N \geq 1} \mathcal{G}_N$ is dense in \mathcal{G} and $\mathcal{G}_N \subset \mathcal{G}_{N+1}$,*
- (c) *Continuity: $E[\log f(p|g; v, r, K)]$ is continuous in g ,*
- (d) *Dominance: There is a function $B(p; v, r, K) > 0$ with $E[B(p; v, r, K)] < \infty$ such that $\log f(p|g; v, r, K) \leq B(p; v, r, K)$ for any g and any $(p; v, r, K)$,*
- (e) *Identification: For any density g with support $(0, \infty)$ such that*

$$E[\log f(p|g; v, r, K)] \geq E[\log f(p|g_0; v, r, K)]$$

$g = g_0$ must hold,

$\lim_{n \rightarrow \infty} \|\hat{g} - g_0\| = 0$ almost surely, provided that $N \equiv N_n \rightarrow \infty$.

This result is a modified version of Theorem 0 in GN. The modification consists of replacing uniform convergence of the objective function by a one-sided uniform convergence implied by Condition (d) and partially by Condition (c), which is possible for maximum likelihood

estimators, as shown by Wald (1949).

In the sequel we discuss briefly how Conditions (a)-(e) can be verified for our problem. Condition (a) follows from Theorem 1 in GN that states that the closure of \mathcal{H} is compact, which is homeomorphic to \mathcal{G} for given γ , σ and by assuming that the location and scale parameters γ , σ are in a compact subset of $\mathbb{R} \times (0, \infty)$. Condition (b) follows from Theorem 2 in GN by using the homeomorphism between \mathcal{H} and \mathcal{G} . Whether Condition (c) is satisfied or not depends on whether the price density $f(p|g; v, r, K)$ is continuous in g . This mild condition appears to be very difficult to verify due to the implicit nature of the price distribution and the nonlinearity of the system of equations that determines the price distribution. Condition (d) is a one-sided dominance condition for which we provide primitive conditions in Lemma A.1 below. These primitive conditions are sufficient for the case when firms' marginal cost r is estimated from an additional data source, so we can regard the valuations and marginal costs in every market as known by the econometrician.¹⁹ Condition (e) can be verified under the conditions of our identification result in Proposition 2 by using the (Shannon-Kolmogorov) Information Inequality.

Lemma A.1 (1) For any density g with support $(0, \infty)$ and any $(p; v, r, K)$

$$\log f(p|g; v, r, K) \leq |\log(v - r)| + |\log(p - r)| + |\log(v - p)| \equiv B(p; v, r, K).$$

(2) Assume that g_0 and the joint distribution of (v, r, K) satisfy the following conditions: (i) $f(v, r, K)$ is bounded; (ii) either (A) g_0 has at least polynomial upper tail, i.e., there is $\alpha > 0$, $L > 0$, $\bar{c} > 1/2$ such that $g_0(c) \geq Lc^{-1-\alpha}$ for any $c > \bar{c}$ and $\int v^{\alpha+1} |\log v| f(v) dv < \infty$, or (B) g_0 has at least exponential upper tail, i.e., there is $\alpha > 0$, $L > 0$, $\bar{c} > 1/2$ such that $g_0(c) \geq Le^{-\alpha c}$ for any $c > \bar{c}$ and the distribution of valuations has at most exponential upper tail, i.e., there is $\alpha' > 0$, $L' > 0$, $\bar{c}' > 1/2$ such that $g_0(c) < L'e^{-\alpha' c}$ for any $c > \bar{c}'$ with $\alpha' > \alpha$.

We note that conditions (i), (ii) are somewhat restrictive, but they still allow the one-sided dominance condition to hold for a large class of search cost and valuation distributions. Condition (ii) suggests that there is a trade-off between the restrictions on the tails of the search

¹⁹In this respect our results are incomplete, but we believe they are still interesting because they serve as an illustration of how one can verify the dominance condition (d) in a structural model so highly nonlinear as ours.

cost and valuation distributions.

Proof of Lemma A.1. For notational simplicity let us drop the conditioning variables v, r, K from $f(p|g; v, r, K)$. From (10)

$$f(p|g) = \frac{\mu_1(v-r)}{(p-r)^2 \sum_{k=2}^K k(k-1) \mu_k (1-F(p|g))^{k-2}}, \quad (\text{A17})$$

and since

$$\sum_{k=2}^K k(k-1) \mu_k (1-F(p|g))^{k-2} \geq \sum_{k=2}^K k \mu_k (1-F(p|g))^{k-1},$$

we obtain

$$\begin{aligned} f(p|g) &\leq \frac{\mu_1(v-r)}{(p-r)^2 \sum_{k=2}^K k \mu_k (1-F(p|g))^{k-1}} \\ &= \frac{\mu_1(v-r)}{(p-r)^2 \sum_{k=1}^K k \mu_k (1-F(p|g))^{k-1} - \mu_1(p-r)^2} \\ &= \frac{\mu_1(v-r)}{\mu_1(p-r)(v-r) - \mu_1(p-r)^2} = \frac{(v-r)}{(p-r)(v-p)}, \end{aligned}$$

where the last-but-one equality follows from (4). That is,

$$\log f(p|g) \leq \log \left[\frac{(v-r)}{(p-r)(v-p)} \right] \leq |\log(v-r)| + |\log(p-r)| + |\log(v-p)| = B(p; v, r, K).$$

This establishes (1).

In what follows we prove (2). We have

$$\begin{aligned} E[B(p; v, r, K)] &= \int (|\log(v-r)| + |\log(p-r)| + |\log(v-p)|) f(p|g_0) f(v, r, K) dp d(v, r, K) \\ &= \int \left[\int_{p_0}^v (|\log(v-r)| + |\log(p-r)| + |\log(v-p)|) f(p|g_0) dp \right] f(v, r, K) d(v, r, K) \\ &\equiv I_1 + I_2 + I_3, \end{aligned} \quad (\text{A18})$$

where $f(v, r, K)$ is the joint density of (v, r, K) . Below we prove that the three integrals I_1, I_2, I_3 are finite.

Bounding I_1 . We have

$$I_1 = \int |\log(v-r)| \left[\int_{p_0}^v f(p|g_0) dp \right] f(v,r,K) d(v,r,K) = \int |\log(v-r)| f(v,r,K) d(v,r,K).$$

This can be split such that

$$\begin{aligned} \int |\log(v-r)| f(v,r,K) d(v,r,K) &= \int_{v-r \leq 1} |\log(v-r)| f(v,r,K) d(v,r,K) \\ &\quad + \int_{v-r > 1} \log(v-r) f(v,r,K) d(v,r,K). \end{aligned}$$

The first term is finite by Condition (i) and the fact that $\int_0^1 |\log x| dx = 1$. The second term is also finite because

$$\begin{aligned} \int_{v-r > 1} \log(v-r) f(v,r,K) d(v,r,K) &< \int_{v > 1} \log(v) f(v,r,K) d(v,r,K) = \int_{v > 1} \log(v) f(v) dv \\ &< \int v f(v) dv, \end{aligned}$$

which is finite by Condition (ii,A). Here and throughout this proof $f(v)$ denotes the marginal density of v .

Bounding I_2 . We have

$$I_2 = \int \left[\int_{p_0}^v |\log(p-r)| f(p|g_0) dp \right] f(v,r,K) d(v,r,K);$$

First focus on the integral in the brackets. Since

$$\sum_{k=2}^K k(k-1) \mu_k (1 - F(p|g))^{k-2} \geq 2\mu_2$$

from (A17) we obtain,

$$f(p|g) \leq \frac{\mu_1(v-r)}{2(p-r)^2 \mu_2} = \left(\frac{v-r}{p-r} \right)^2 f(p|g)|_{p=v}. \quad (\text{A19})$$

Then

$$\begin{aligned} \int_{\underline{p}_0}^v |\log(p-r)| f(p|g_0) dp &\leq \int_{\underline{p}_0}^v |\log(p-r)| \left(\frac{v-r}{p-r}\right)^2 f(p|g_0)|_{p=v} dp \\ &= (v-r) f(p|g_0)|_{p=v} \int_{\underline{p}_0}^v |\log(p-r)| \frac{v-r}{(p-r)^2} dp, \end{aligned}$$

where

$$\begin{aligned} \int_{\underline{p}_0}^v |\log(p-r)| \frac{v-r}{(p-r)^2} dp &= \int_1^{\frac{v-r}{\underline{p}_0-r}} \left| \log\left(\frac{v-r}{x}\right) \right| dx \\ &\leq \int_1^{\frac{v-r}{\underline{p}_0-r}} |\log(v-r)| dx + \int_1^{\frac{v-r}{\underline{p}_0-r}} \log x dx \\ &= |\log(v-r)| \left(\frac{v-r}{\underline{p}_0-r} - 1\right) + \frac{v-r}{\underline{p}_0-r} \left(\log \frac{v-r}{\underline{p}_0-r} - 1\right) + 1 \\ &\leq |\log(v-r)| \frac{v-r}{\underline{p}_0-r} + \frac{v-r}{\underline{p}_0-r} \log \frac{v-r}{\underline{p}_0-r} + 1. \end{aligned} \quad (\text{A20})$$

So

$$\int_{\underline{p}_0}^v |\log(p-r)| f(p|g_0) dp \leq (v-r) f(p|g_0)|_{p=v} \left[|\log(v-r)| \frac{v-r}{\underline{p}_0-r} + \frac{v-r}{\underline{p}_0-r} \log \frac{v-r}{\underline{p}_0-r} + 1 \right].$$

Based on this, we need to show that

$$J_1 = \int (v-r) f(p|g_0)|_{p=v} |\log(v-r)| \frac{v-r}{\underline{p}_0-r} f(v,r,K) d(v,r,K) < \infty, \quad (\text{A21})$$

$$J_2 = \int (v-r) f(p|g_0)|_{p=v} \frac{v-r}{\underline{p}_0-r} \log \frac{v-r}{\underline{p}_0-r} f(v,r,K) d(v,r,K) < \infty, \quad (\text{A22})$$

$$J_3 = \int (v-r) f(p|g_0)|_{p=v} f(v,r,K) d(v,r,K) < \infty. \quad (\text{A23})$$

We expect that $f(p|g_0)|_{p=v} < M$ for some appropriate M for any (v,r,K) because $f(p|g_0)|_{p=v} \rightarrow 0$ when $v \rightarrow \infty$, since $f(p|g_0)|_{p=v}$ is the density at the upper bound of its support, although we find it difficult to prove this formally. Further, by (5)

$$\frac{v-r}{\underline{p}_0-r} = \frac{\sum_{k=1}^K k \mu_{k0}}{\mu_{10}} > 1,$$

where $(\mu_{k0})_{k=1}^K$ correspond to the true g_0 . The numerator is bounded, in fact $\sum_{k=1}^K k\mu_{k0} \in [1, K]$ for any g_0 . By (3a),

$$\frac{1}{\mu_{10}} = \frac{1}{1 - G_0(c_{10})} \leq \frac{1}{1 - G_0\left(\frac{v-r}{2}\right)}$$

because G_0 is increasing and from (8)

$$\begin{aligned} c_{10} &= \int_0^1 \left(\frac{\mu_1(v-r)}{\sum_{k=1}^K k\mu_k(1-z)^{k-1}} + r \right) (2z-1) dz \leq \mu_1(v-r) \int_0^1 \frac{|2z-1|}{\sum_{k=1}^K k\mu_k(1-z)^{k-1}} dz \\ &\leq (v-r) \int_0^1 |2z-1| dz = \frac{v-r}{2}, \end{aligned}$$

where the latter inequality follows by taking $z = 1$ in the denominator. Therefore,

$$\frac{v-r}{\underline{p}_0 - r} \leq \frac{K}{1 - G_0\left(\frac{v-r}{2}\right)}. \quad (\text{A24})$$

Now we proceed by proving (A21)-(A23). Applying (A24), we have

$$\begin{aligned} J_1 &\leq MK \int \frac{(v-r) |\log(v-r)|}{1 - G_0\left(\frac{v-r}{2}\right)} f(v, r, K) d(v, r, K) \\ &= MK \int_{v-r \leq 2\bar{c}} \frac{(v-r) |\log(v-r)|}{1 - G_0\left(\frac{v-r}{2}\right)} f(v, r, K) d(v, r, K) \\ &\quad + MK \int_{v-r > 2\bar{c}} \frac{(v-r) \log(v-r)}{1 - G_0\left(\frac{v-r}{2}\right)} f(v, r, K) d(v, r, K). \end{aligned}$$

The first term is finite because the function $x \log x$ is bounded on any bounded interval and $1/[1 - G_0\left(\frac{v-r}{2}\right)] \leq 1/[1 - G_0(\bar{c})]$. For the second term we note that $x \log x / [1 - G_0\left(\frac{x}{2}\right)]$ is an increasing function in x , so

$$\begin{aligned} \int_{v-r > 2\bar{c}} \frac{(v-r) \log(v-r)}{1 - G_0\left(\frac{v-r}{2}\right)} f(v, r, K) d(v, r, K) &< \int_{v > 2\bar{c}} \frac{v \log v}{1 - G_0\left(\frac{v}{2}\right)} f(v, r, K) d(v, r, K) \\ &= \int_{v > 2\bar{c}} \frac{v \log v}{1 - G_0\left(\frac{v}{2}\right)} f(v) dv. \end{aligned}$$

Under Condition (ii,A), for $c > \bar{c}$

$$1 - G_0(c) = \int_c^\infty g_0(x) dx \geq \int_c^\infty Lx^{-1-\alpha} dx = L \frac{c^{-\alpha}}{\alpha},$$

so for $v > 2\bar{c}$

$$1 - G_0\left(\frac{v}{2}\right) \geq 2^\alpha L \frac{v^{-\alpha}}{\alpha}. \quad (\text{A25})$$

Therefore the second term of J_2 is less than

$$\alpha 2^{-\alpha} L^{-1} MK \int_{v>2\bar{c}} v^{\alpha+1} \log(v) f(v) dv < \alpha 2^{-\alpha} L^{-1} MK \int v^{\alpha+1} |\log v| f(v) dv < \infty,$$

the latter inequality by the second part of Condition (ii,A). This proves $J_1 < \infty$.

Under Condition (ii,B), for $c > \bar{c}$

$$1 - G_0(c) = \int_c^\infty g_0(x) dx \geq \int_c^\infty L e^{-\alpha x} dx = L \frac{e^{-\alpha c}}{\alpha},$$

so for $v > 2\bar{c}$

$$1 - G_0\left(\frac{v}{2}\right) \geq \frac{L}{\alpha} e^{-\frac{\alpha v}{2}}. \quad (\text{A26})$$

Therefore the second term of J_2 is less than

$$\alpha L^{-1} MK \int_{v>2\bar{c}} v \log(v) e^{\frac{\alpha v}{2}} f(v) dv < \alpha L^{-1} L' MK \int_{v>2\bar{c}} v \log(v) e^{-(\alpha' - \frac{\alpha}{2})v} dv < \infty,$$

where the former inequality follows from the second part of Condition (ii,B). This proves $J_1 < \infty$.

Now, applying again (A24), we have

$$J_2 \leq MK \int \frac{(v-r)}{1 - G_0\left(\frac{v-r}{2}\right)} \log\left(\frac{K}{1 - G_0\left(\frac{v-r}{2}\right)}\right) f(v, r, K) d(v, r, K). \quad (\text{A27})$$

This can be split into the sum of two integrals:

$$\begin{aligned} & MK \int_{v-r \leq 2\bar{c}} \frac{(v-r)}{1 - G_0\left(\frac{v-r}{2}\right)} \log\left(\frac{K}{1 - G_0\left(\frac{v-r}{2}\right)}\right) f(v, r, K) d(v, r, K) \\ & + MK \int_{v-r > 2\bar{c}} \frac{(v-r)}{1 - G_0\left(\frac{v-r}{2}\right)} \log\left(\frac{K}{1 - G_0\left(\frac{v-r}{2}\right)}\right) f(v, r, K) d(v, r, K). \end{aligned}$$

The first integral is less than

$$MK \frac{2\bar{c}}{1 - G_0(\bar{c})} \log \left(\frac{K}{1 - G_0(\bar{c})} \right) < \infty.$$

The second integral can be bounded in a way similar to the second integral term of J_1 . We obtain

$$\begin{aligned} & \int_{v-r > 2\bar{c}} \frac{(v-r)}{1 - G_0\left(\frac{v-r}{2}\right)} \log \left(\frac{K}{1 - G_0\left(\frac{v-r}{2}\right)} \right) f(v, r, K) d(v, r, K) \\ & < \int_{v > 2\bar{c}} \frac{v}{1 - G_0\left(\frac{v}{2}\right)} \log \left(\frac{K}{1 - G_0\left(\frac{v}{2}\right)} \right) f(v, r, K) d(v, r, K) \\ & = \int_{v > 2\bar{c}} \frac{v}{1 - G_0\left(\frac{v}{2}\right)} \log \left(\frac{K}{1 - G_0\left(\frac{v}{2}\right)} \right) f(v) dv. \end{aligned} \quad (\text{A28})$$

Under Condition (ii,A), from (A25) this is less than

$$\begin{aligned} & \alpha 2^{-\alpha} L^{-1} \int_{v > 2\bar{c}} v^{\alpha+1} \log \left(\frac{\alpha K}{2^\alpha L} v^\alpha \right) f(v) dv \\ & = \alpha 2^{-\alpha} L^{-1} \int_{v > 2\bar{c}} v^{\alpha+1} (\log a + \alpha \log v) f(v) dv \\ & = \alpha 2^{-\alpha} L^{-1} \log a \int_{v > 2\bar{c}} v^{\alpha+1} f(v) dv + \alpha^2 2^{-\alpha} L^{-1} \int_{v > 2\bar{c}} v^{\alpha+1} |\log v| f(v) dv \\ & < (\alpha 2^{-\alpha} L^{-1} \log a + \alpha^2 2^{-\alpha} L^{-1}) \int_{v > 2\bar{c}} v^{\alpha+1} |\log v| f(v) dv, \end{aligned}$$

where $a = \alpha 2^{-\alpha} K L^{-1}$. Consequently, Condition (ii,A) implies that this is finite, and therefore $J_2 < \infty$.

Under Condition (ii,B), from (A26) the expression in (A28) is less than

$$\begin{aligned} & \alpha L^{-1} \int_{v > 2\bar{c}} v e^{\frac{\alpha v}{2}} \log \left(\frac{\alpha K}{L} e^{\frac{\alpha v}{2}} \right) f(v) dv \\ & = \alpha L^{-1} \int_{v > 2\bar{c}} v e^{\frac{\alpha v}{2}} \left(\log a + \frac{\alpha v}{2} \right) f(v) dv \\ & = \alpha L^{-1} \log a \int_{v > 2\bar{c}} v e^{\frac{\alpha v}{2}} f(v) dv + \frac{\alpha^2 L^{-1}}{2} \int_{v > 2\bar{c}} v^2 e^{\frac{\alpha v}{2}} f(v) dv \\ & < \alpha L^{-1} L' \log a \int_{v > 2\bar{c}} v e^{-(\alpha' - \frac{\alpha}{2})v} dv + \frac{\alpha^2 L^{-1} L'}{2} \int_{v > 2\bar{c}} v^2 e^{-(\alpha' - \frac{\alpha}{2})v} dv < \infty, \end{aligned}$$

where $a = \alpha KL^{-1}$. Consequently, $J_2 < \infty$.

The statement in (A23) follows easily from the second part of Condition (ii,A). This completes the proof of $I_2 < \infty$.

Bounding I_3 . We have

$$I_3 = \int \left[\int_{\underline{p}_0}^v |\log(v-p)| f(p|g_0) dp \right] f(v,r,K) d(v,r,K).$$

The integral in the brackets is

$$\begin{aligned} \int_{\underline{p}_0}^v |\log(v-p)| f(p|g_0) dp &\leq \int_{\underline{p}_0}^v |\log(v-p)| \left(\frac{v-r}{p-r} \right)^2 f(p|g_0) |_{p=v} dp \\ &= (v-r) f(p|g_0) |_{p=v} \int_{\underline{p}_0}^v |\log(v-p)| \frac{v-r}{(p-r)^2} dp, \end{aligned}$$

where

$$\begin{aligned} \int_{\underline{p}_0}^v |\log(v-p)| \frac{v-r}{(p-r)^2} dp &= \int_1^{\frac{v-r}{\underline{p}_0-r}} \left| \log \left[(v-r) \frac{x-1}{x} \right] \right| dx \\ &\leq \int_1^{\frac{v-r}{\underline{p}_0-r}} |\log(v-r)| dx + \int_1^{\frac{v-r}{\underline{p}_0-r}} \log x dx - \int_1^{\frac{v-r}{\underline{p}_0-r}} \log(x-1) dx \\ &= |\log(v-r)| \left(\frac{v-r}{\underline{p}_0-r} - 1 \right) + \frac{v-r}{\underline{p}_0-r} \log \frac{v-r}{\underline{p}_0-r} \\ &\quad - \left(\frac{v-r}{\underline{p}_0-r} - 1 \right) \log \left(\frac{v-r}{\underline{p}_0-r} - 1 \right) \\ &\leq |\log(v-r)| \frac{v-r}{\underline{p}_0-r} + \frac{v-r}{\underline{p}_0-r} \log \frac{v-r}{\underline{p}_0-r} \\ &\quad - \left(\frac{v-r}{\underline{p}_0-r} - 1 \right) \log \left(\frac{v-r}{\underline{p}_0-r} - 1 \right). \end{aligned} \tag{A29}$$

So we need to show that

$$\begin{aligned} H_1 &= \int (v-r) f(p|g_0) |_{p=v} |\log(v-r)| \frac{v-r}{\underline{p}_0-r} f(v,r,K) d(v,r,K) < \infty, \\ H_2 &= \int (v-r) f(p|g_0) |_{p=v} \left[\frac{v-r}{\underline{p}_0-r} \log \frac{v-r}{\underline{p}_0-r} \right. \end{aligned} \tag{A30}$$

$$- \left(\frac{v-r}{\underline{p}_0-r} - 1 \right) \log \left(\frac{v-r}{\underline{p}_0-r} - 1 \right) \Big] f(v, r, K) d(v, r, K) < \infty. \quad (\text{A31})$$

The first statement is proved in (A21). For the second statement we note that the function $x \log x - (x-1) \log(x-1)$ is increasing in x . Therefore, by (A24)

$$\begin{aligned} H_2 &< M \int (v-r) \left[\frac{K}{1-G_0\left(\frac{v-r}{2}\right)} \log \frac{K}{1-G_0\left(\frac{v-r}{2}\right)} \right. \\ &\quad \left. - \left(\frac{K}{1-G_0\left(\frac{v-r}{2}\right)} - 1 \right) \log \left(\frac{K}{1-G_0\left(\frac{v-r}{2}\right)} - 1 \right) \right] f(v, r, K) d(v, r, K) \\ &< MK \int \frac{(v-r)}{1-G_0\left(\frac{v-r}{2}\right)} \log \frac{K}{1-G_0\left(\frac{v-r}{2}\right)} f(v, r, K) d(v, r, K). \end{aligned}$$

The latter expression is the same as the RHS expression in (A27), which we have already proved to be finite. Consequently, $I_3 < \infty$. This completes the proof that $E[B(p; v, r, K)] < \infty$. ■

Figures

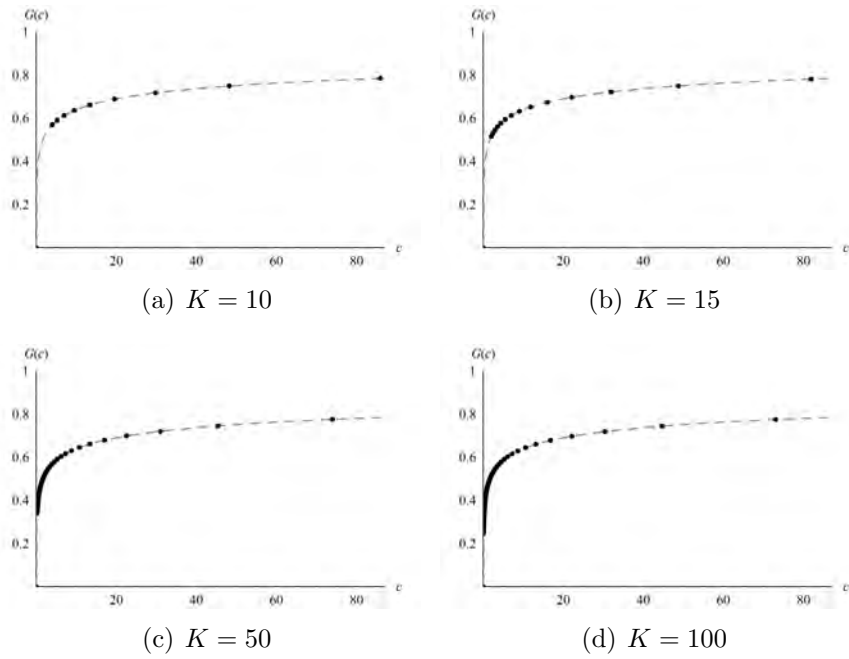


Figure 1: Search cost cutoffs with data from only one market

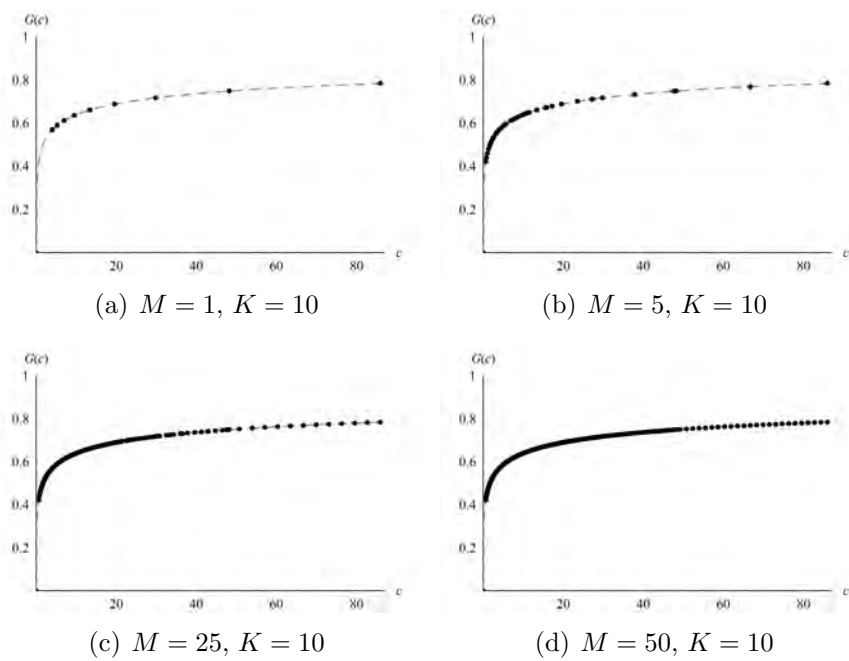


Figure 2: Search cost cutoffs with data from M different markets

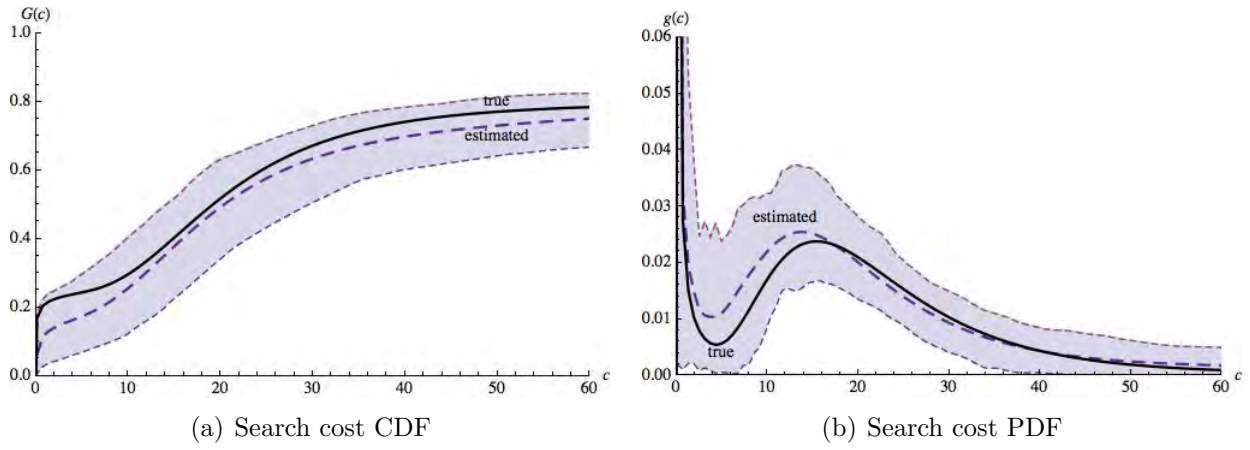


Figure 3: Monte Carlo results: estimated search costs for $N = 8$)

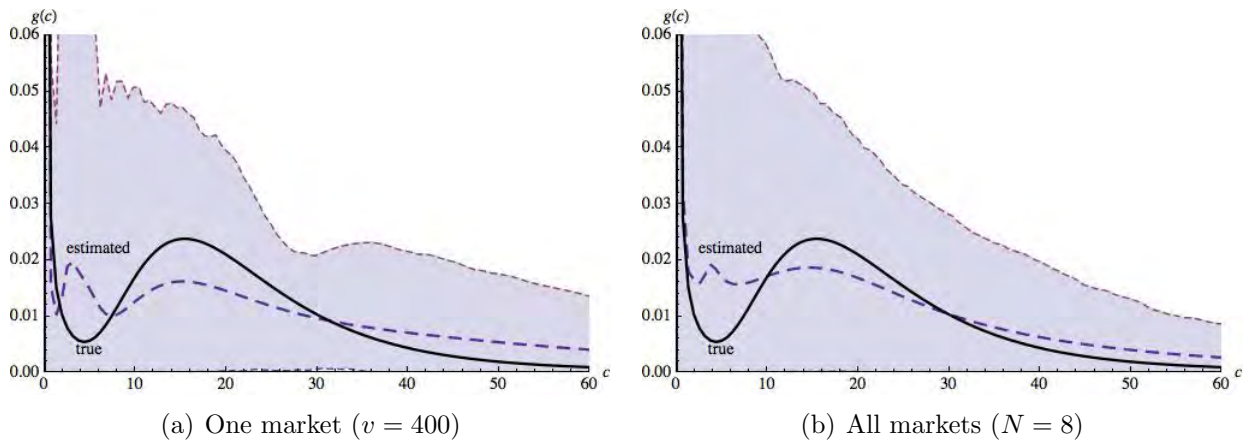


Figure 4: Monte Carlo results: estimated search costs (market-by-market)

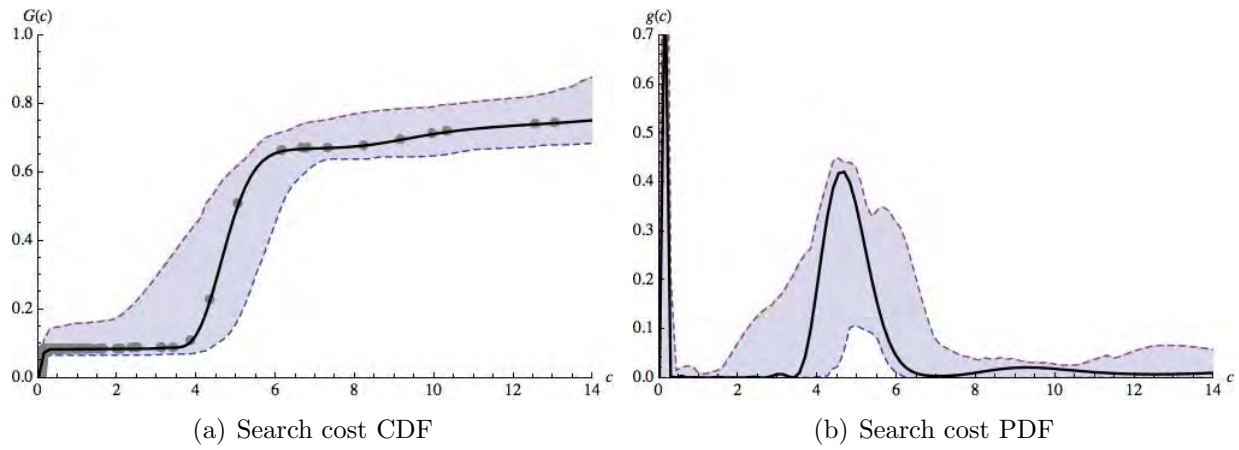


Figure 5: Estimated search cost distribution

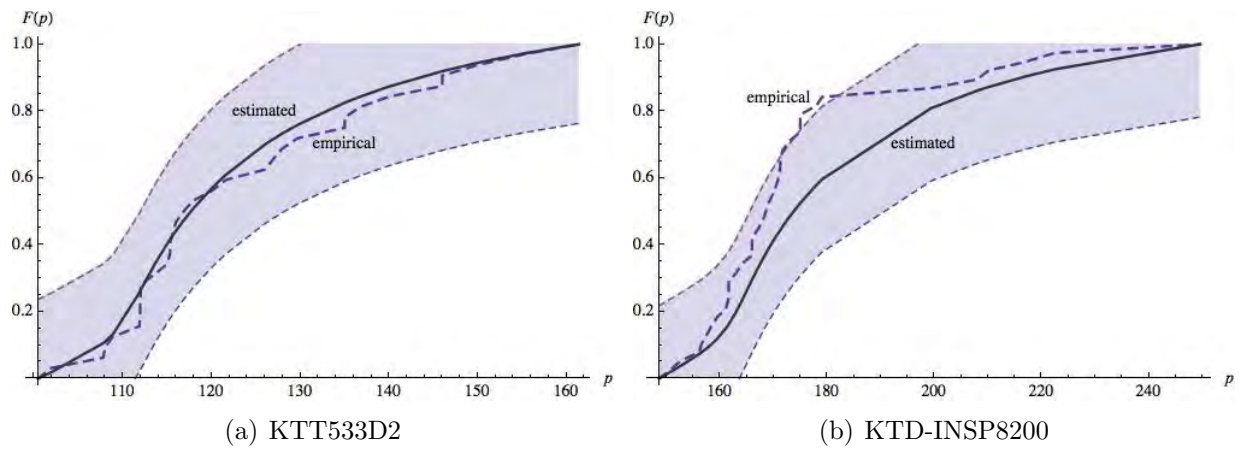
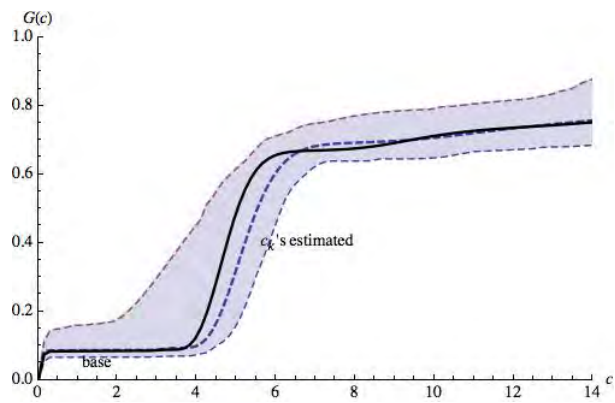
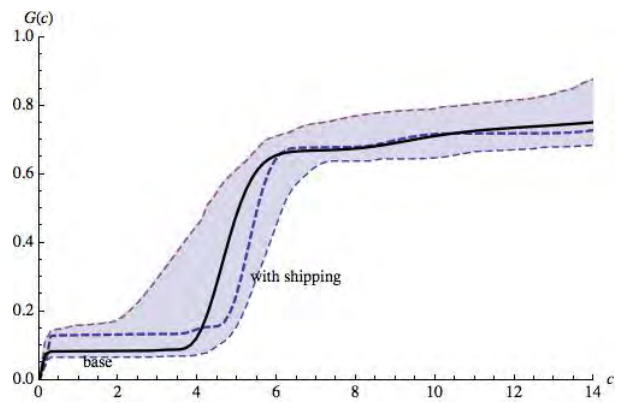


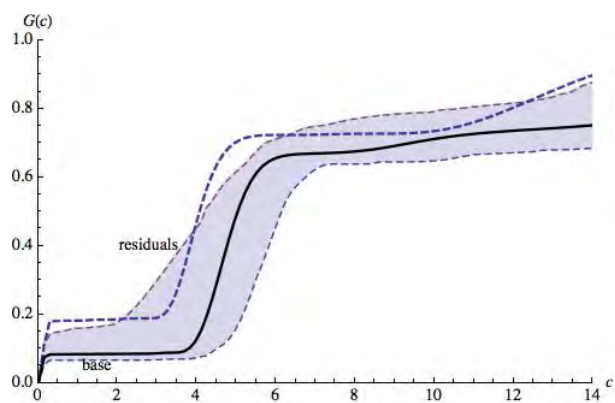
Figure 6: Estimated and empirical price CDF



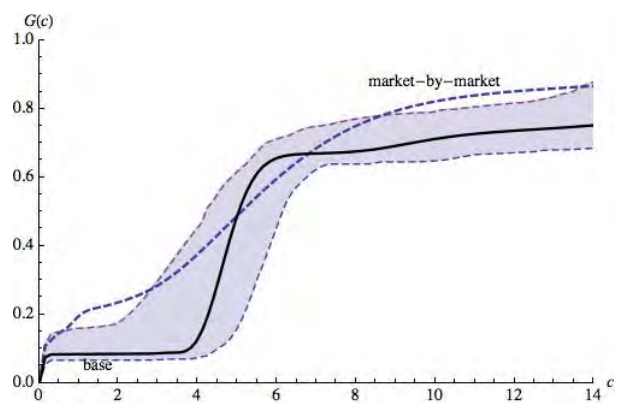
(a) Search cost cutoffs estimated



(b) Including shipping costs



(c) Residuals



(d) Market-by-market approach

Figure 7: Estimated search cost CDF alternative specifications

Tables

Table 1: Monte Carlo results

	Prices		Search costs
	ISE (approx.)	ISE (true) $\times 10^{-2}$	ISE (true) $\times 10^{-4}$
$N = 1$	-0.097	1.031	1.625
$N = 2$	-0.235	0.838	0.181
$N = 3$	-0.236	0.841	0.147
$N = 4$	-0.244	0.844	0.143
$N = 5$	-0.241	0.829	0.121
$N = 6$	-0.242	0.822	0.097
$N = 7$	-0.243	0.824	0.092
$N = 8$	-0.245	0.823	0.077
$N = 9$	-0.244	0.839	0.078

Notes: ISE values are calculated for the mean price and search cost densities of the 100 replications.

Table 2: List of products

Part number	Manufacturer	Compatibility	Size	Speed	Form factor
KTT3311A	Kingston	Toshiba	1GB	333MHz DDR333/PC2700	200-pin SoDIMM
KTT533D2	Kingston	Toshiba	1GB	533MHz DDR2-533/PC2-4200	200-pin SoDIMM
KTD-INSP8200	Kingston	Dell	1GB	266MHz DDR266/PC2100	200-pin SoDIMM
KTD-INSP5150	Kingston	Dell	1GB	333MHz DDR333/PC2700	200-pin SoDIMM
KTD-INSP6000	Kingston	Dell	1GB	533MHz DDR2-533/PC2-4200	240-pin SoDIMM
KTD-INSP6000A	Kingston	Dell	1GB	533MHz DDR2-533/PC2-4200	200-pin SoDIMM
KAC-MEME	Kingston	Acer	1GB	533MHz DDR2-533/PC2-4200	200-pin SoDIMM
KTD-INSP9100	Kingston	Dell	1GB	400MHz DDR400/PC3200	200-pin SoDIMM
KTM-TP3840	Kingston	IBM	1GB	533MHz DDR2-533/PC2-4200	200-pin SoDIMM
KTH-ZD8000A	Kingston	HP Compaq	1GB	533MHz DDR2-533/PC2-4200	200-pin SoDIMM

Table 3: Summary statistics

Part number	No. of Stores	Mean Price (Std)	Min. Price	Max. Price	Coeff. of Var. (as %)
KTT3311A	32	181.67 (24.62)	148.62	235.00	13.55
KTT533D2	33	123.33 (15.62)	100.45	161.40	12.66
KTD-INSP8200	39	173.59 (21.31)	148.62	249.54	12.28
KTD-INSP5150	39	179.09 (19.84)	148.62	222.35	11.08
KTD-INSP6000	35	120.29 (13.48)	100.45	151.05	11.21
KTD-INSP6000A	38	116.33 (13.43)	94.99	154.50	11.54
KAC-MEME	24	123.58 (17.47)	101.92	161.64	14.14
KTD-INSP9100	33	175.84 (24.38)	148.62	249.54	13.87
KTM-TP3840	37	122.83 (14.32)	104.55	161.94	11.65
KTH-ZD8000A	41	116.77 (12.25)	100.45	154.50	10.49

Notes: Prices are in US dollars.

Table 4: Fit SNP estimates for different values of N

N	LL	ISE	KS (avg)
1	40.047	-0.0416	1.289
2	40.011	-0.0414	1.286
3	39.532	-0.0449	1.085
4	39.527	-0.0451	1.089
5	39.525	-0.0451	1.091
10	39.258	-0.0484	1.073
15	39.207	-0.0487	1.038
20	39.180	-0.0492	1.044
25	39.158	-0.0491	1.038
30	39.146	-0.0492	1.035
35	39.141	-0.0493	1.046
40	39.132	-0.0492	1.029
45	39.131	-0.0493	1.041
50	39.127	-0.0493	1.029

Notes: c_k 's are obtained from empirical price CDF.

Table 5: Parameter estimates products and fit

Part number	K	p	v	r	μ_1	μ_2	μ_3	μ_4	$\mu_{5...15}$	$\mu_{16...K}$	KS
KTT3311A	32	148.62	235.00	142.73	0.26	0.24	0.42	0.00	0.00	0.08	0.66
KTT533D2	33	100.45	161.40	93.93	0.32	0.59	0.00	0.00	0.00	0.08	0.61
KTD-INSP8200	39	148.62	249.54	138.75	0.29	0.62	0.00	0.00	0.00	0.08	1.69
KTD-INSP5150G	39	148.62	222.35	142.52	0.28	0.61	0.02	0.00	0.00	0.08	1.63
KTD-INSP6000	35	100.45	151.05	94.74	0.33	0.58	0.00	0.00	0.00	0.08	0.95
KTD-INSP6000A	38	94.99	154.50	87.40	0.33	0.58	0.00	0.00	0.01	0.08	1.16
KAC-MEME	24	101.92	161.64	96.28	0.26	0.61	0.00	0.00	0.01	0.07	0.73
KTD-INSP9100	33	148.62	249.54	139.11	0.26	0.51	0.14	0.00	0.04	0.04	1.03
KTM-TP3840	37	104.55	161.94	97.09	0.33	0.59	0.00	0.00	0.02	0.06	0.84
KTH-ZD8000A	41	100.45	154.50	93.30	0.34	0.58	0.00	0.00	0.02	0.07	1.15