# Estimating Store Choices with Endogenous Shopping Bundles and Price Uncertainty

Hyunchul Kim*　　　　　　　Kyoo il Kim

Sungkyunkwan University　　Michigan State University

July 2015

### Abstract

We develop a structural model of retail store choices for which household shopping plans and price beliefs are endogenously determined. In our model individual households make their store choices based on their expected basket costs, which are determined by their shopping plans and price beliefs. Previous studies use realized purchases as a proxy for unobserved shopping lists and also assume homogenous price expectation across all households over the entire sample period. Our approach improves the measures of expected basket costs by estimating intended shopping lists of households using a duration model and also by constructing household-, time-, store-, and goods-specific price expectations. In our empirical application using a scanner data set, we find that the store choices become significantly more elastic to prices when the correction is applied.

Keywords: Store choice, Expected basket cost, Price expectation, Consumer demand, Supermarket industry

JEL Classification: D1, D8, L1, L8

# 1 Introduction

In this paper we develop a household level demand model of store choice in the supermarket industry. The store choice of grocery shoppers has two distinctive features compared to other consumer choice problems. First, store choice decisions involve bundle purchase behavior. Consumers decide on which stores to visit depending on their shopping plans, which are characterized by the goods and quantities they intend to buy at a store. Second, consumers face price uncertainty before they actually visit a store. Shoppers may acquire certain information on prices from out-of-store advertising such as newspaper inserts or weekly circulars on sales, but most of the shelf prices are unknown a priori. Shoppers, therefore, may rely on price expectation for their store choice from their shopping experiences.

Although household level *planned* shopping lists and price beliefs are central to store choice decisions, these are typically not observed by researchers. Previous store choice studies used realized purchases as a proxy for the unobserved shopping lists and assumed homogeneous price expectations across households and over the entire sample period (see Bell, Ho, and Tang (1998), Smith (2004), Beresteanu, Ellickson, and Misra (2010), Briesch, Chintagunta, and Fox (2009), and Hansen and Singh (2009)). This approach introduces potential measurement problems in the expected basket costs of each store trip and thus may create biases in elasticity estimates of expected basket costs for store choices. The main contribution of this study is to provide an empirical method of estimating household level store choices using improved measures of the expected basket costs. To do this, we model household level shopping baskets and also construct household-, store-, and time-specific price expectations based on the past shopping experiences of each household. To our knowledge no prior work examines the potential biases caused by unobserved household shopping lists and price beliefs in the estimation of store choice.

To improve the measure of unobserved shopping lists, we proceed in the following steps. In the first stage, we estimate a model of how households determine which goods to buy and how much to buy before they visit a store.[1] In particular we use a continuous-time duration

---

[1]Shopping lists may be defined at a narrower level such as brand or brand-size combination. This approach would be more tenable if strong brand loyalty prevails in grocery shopping because in that case shoppers plan on buying specific products. However, brand loyalty is only weakly present in our data for most goods.

approach to model the purchase incidence for each good, which is specified as a function of the relevant states of each household such as consumption and inventory levels and shopping behavior. The quantity choice in planning a shopping list is then predicted as a function of expected prices and other relevant variables. In the second stage, we estimate a household level store choice using the shopping baskets predicted in the first stage as an input. Expected basket costs based on the expected shopping baskets and other store attributes (e.g., transportation costs, service quality, floor space, and parking space) enter a utility function for each store visit and households choose a store that would generate the highest utility from shopping at the store for the expected shopping basket.

The consideration of potential discrepancies between observed and planned shopping baskets in store choice studies was first introduced by Bell, Ho, and Tang (1998). In their seminal study of household level store choice, they assume that households specify shopping lists at the product level assuming strong brand loyalty, as opposed to at the good category level.[2] They use a discrete choice setting to estimate the probability that each product actually purchased was included in the *ex ante* shopping list as a function of consumption, inventory level, and expectation errors in prices. In their purchase incidence estimation, consumption and inventory levels (and the coefficient variants across different customer segments) are the only source of heterogeneity because expected price in their model is a simple average of store prices over the entire sample period and thus does not capture the differences in heterogenous price beliefs across households. In our study, we instead develop the basket composition model at the goods level, which allows for heterogeneity in inter-purchase time for each good as well as consumption and inventory levels.

Given that the idiosyncrasy in purchasing time patterns across households and across goods is salient in grocery shopping, using a duration model that takes into account heterogeneity in inter-purchase time seems more appropriate than a standard discrete choice model in our setting. Moreover, the basket composition model in this study allows for quantity choices for the basket goods, whereas existing store choice studies assume that expected quantities are identical to

---

[2]For example, a shopping list might consist of 144 oz Pepsi cans and 12 oz Kellogg's Special K, instead of soft-drinks and cereal goods.

realized quantities.

Modeling shopping lists not only mitigates the measurement problem, but also importantly it allows us to circumvent the complication of estimating store choices in a discrete choice setup, which arises from the size and complexity of choice sets (see Katz (2007) and Pakes (2010)). The choice set in a store choice problem not merely consists of alternative stores, but it also contains alternative shopping baskets (what and how much to buy) depending on which store is chosen. Therefore, the choice set for each store trip should include all the combinations of alternative stores and the shopping baskets corresponding to each store, which would make the store choice estimation in a discrete choice setting almost intractable. Specifying the expected shopping baskets in the first stage allows us to sidestep the need to deal with such overwhelming size and complexity of the choice sets.

Next, to deal with unobserved prices, we develop household level price expectations for each good as a function of most recent trips to the store. It is well established in the marketing and economics literature that households glean price information from external advertising or by retrieving memories of store prices (see, for example, Kalwani and Yim (1992), Erdem, Imai, and Keane (2003), and Hendel and Nevo (2006)). We assume that price expectation builds on store prices experienced from past shopping trips and purchases. Price expectation constructed this way not only captures short-term fluctuations in store prices but also reflects heterogeneity in price beliefs across households depending on their shopping behavior. Bell, Ho, and Tang (1998) consider price uncertainty and assume that customers have information on the price distribution, particularly the average of store prices (see also Alba, Broniarczyk, Shimp, and Urbany (1994), Lal and Rao (1997), and Ho, Tang, and Bell (1998)). This approach is grounded on the rational expectation which imposes homogeneity in price beliefs. Therefore, this approach of specifying price expectation – that are common across different households – does not capture heterogeneity in price beliefs among households.

Using our proposed method we estimate a store choice model using a scanner data set. We find that when our approach is applied to household price expectation and planned shopping lists, the store-level own price elasticities become dramatically higher. Particularly, the potential measurement problems in expected basket costs are most attributable to ignoring heterogeneity

4

in price expectation across households. Estimating the shopping list at the product level, as opposed to at the goods level, somewhat reduces such biases in the price elasticity estimates even without heterogeneity in price beliefs, but the own price elasticities on average are less than half those based on our approach in magnitude.

Correctly estimating consumer demand and price elasticities is central to the study of many important problems such as pricing strategies of firms, antitrust policy, and welfare effects of introducing new goods. For instance, using the estimates of demand parameters store or brand managers, who design pricing policies with short-term price promotion, can predict how much price cuts increase future store traffics or accelerate in-store purchases. Estimating consumer demand and price elasticities also plays a key role in understanding the welfare implication of merger policies or zoning regulations that restrict new stores' entry by providing precise measures of the extent of competition between retailers. This is because the demand parameter determines consumer choice and firm profitability under counterfactual industry structure. Therefore understanding the demand system is the first step toward the analysis of firms incentives to enter/exit as well as their location choices.

The rest of the paper is organized as follows. Section 2 outlines our store choice model. In Section 3, we develop a model of the expected basket costs and provide estimation methods. Section 4 describes the scanner data of consumer choices. The estimation results are presented in Section 5. Section 6 concludes.

## 2    Store Choice Model

We use a standard characteristics-based random utility framework to model household level store choices. Following the characteristics based approach (Lancaster (1971) and McFadden (1981)), the demand system of store choice posits that choice patterns are determined by preferences for store attributes. We specify the utility of household $i$ that plans to shop for basket $b$ from

visiting store $j$ at time $t$ as[3]

$$U_{ijbt} = \beta_i Dist_{ij} - \alpha_i Exp_{ijbt} + \gamma_{ij} + \varepsilon_{ijt}, \tag{1}$$

where $Dist_{ij}$ is the shortest driving distance between store $j$ and the residence of household $i$, $Exp_{ijbt}$ is expected spending on shopping bundle $b$, and $\gamma_{ij}$ is the marginal utility to unobservable (to the econometrician) attributes of store $j$. $\varepsilon_{ijt}$ is the random part of the utility that captures the idiosyncratic taste of household $i$ for store $j$.

Here the expected bundle cost $Exp_{ijbt}$ is an important ingredient of the store choice, which is not directly observable from data. Previous studies on store choices (see e.g. Bell, Ho, and Tang (1998), Alba, Broniarczyk, Shimp, and Urbany (1994), Lal and Rao (1997), and Ho, Tang, and Bell (1998)) took several different approaches to approximate this expected bundle cost term using observed data. In this paper we want to improve this crucial measure by taking several novel econometric approaches to the components of the expected bundle cost. We elaborate on this in the following section.

In the utility $\alpha_i$, $\beta_i$, and $\gamma_{ij}$ are household specific parameters to reflect household heterogeneity. We assume these parameters are normally distributed conditional on demographic characteristics as

$$\begin{pmatrix} \alpha_i \\ \beta_i \\ \gamma_{ij} \end{pmatrix} = \begin{pmatrix} \alpha_0 \\ \beta_0 \\ \gamma_{0j} \end{pmatrix} + \Lambda Z_i + \nu_i, \qquad \nu_i \sim N(0, \Sigma) \tag{2}$$

where $Z_i$ is a $d_z \times 1$ vector of demographic characteristics of the household, $\Lambda$ is a $k \times d_z$ matrix of parameters that vary by demographics where $k$ denotes the dimension of the vector of observable characteristics (or store dummies), and $\Sigma$ is a variance-covariance matrix of the multivariate normal distribution. For example, when $Z_i$ includes dummies for income groups, the average marginal utility to income is $\alpha_0$ for household in the base income group and $\alpha_0 + \alpha_g$ for those in the income group $g$.[4] Using this heterogeneous parameter setting with random

---

[3]For a robustness check, we discuss an alternative specification with other basket-related variables. Particularly, we add preferences for a variety of products in each basket good. See Appendix A.

[4]Using the preferences that vary by demographics, to allow for unobserved store attributes, is in line with Berry, Levinsohn, and Pakes (2004) and Goolsbee and Petrin (2004). In their models, the price coefficient for the base demographic group is subsumed in the product fixed effect term and the estimation requires a two-step

coefficients, our model specification can accommodate important substitution patterns in store choice by allowing consumer heterogeneity in preferences for store attributes.

Last, the idiosyncratic error $\varepsilon_{ijt}$ in equation (1) accounts for the preference shock not captured by observed or unobserved store attributes and demographics. From the random utility (1), we obtain the model choice probability that household $i$ with planned shopping list $b$ chooses store $j$ (at time $t$)

$$s_{ijbt} = \int \left\{ (\nu_i, \varepsilon_{it}) : U_{ijbt} > U_{ij'bt}, \ \forall j' \neq j \right\} dF_\nu(\nu_i) dF_\varepsilon(\varepsilon_{it}), \tag{3}$$

where $F_\nu(\cdot)$ and $F_\varepsilon(\cdot)$ denote the distribution functions of $\nu_i$ and $\varepsilon_{it} = (\varepsilon_{i0t}, \ldots, \varepsilon_{iJt})$, respectively, and $U_{i0bt} = \varepsilon_{i0t}$ denotes the utility of the outside choice (i.e. not going to any store among $j = 1, \ldots, J$). The error term $\varepsilon_{ijt}$ is assumed to be an $i.i.d.$ error with a Type I extreme value distribution. This assumption reduces the model choice probability to a logit model that allows a closed-form solution for the integration over $\varepsilon_{it}$.[5]

# 3    Modeling Expected Bundle Cost

Expected basket cost entering into the utility function in the store choice model is determined by an *ex ante* (or planned) shopping list and household price expectation for the goods included in the shopping list. Since neither of these components is directly observable from the data, estimation of store choice problem should be preceded by specifying *ex ante* shopping lists and price expectations. An *ex ante* shopping list is characterized by the set of goods and the quantities of the goods that a shopper plans to purchase before visiting a store. The choice of a shopping list is specified in two stages. In the first stage, the set of goods included in a shopping list is determined as a function of inter-purchase spells and other relevant states such as inventory and consumption level for each good. In the second stage, given this set of goods

---

procedure to identify price parameters. In contrast, in our setting, the variation in the expected basket costs across households and shopping baskets allows the model to identify price parameters for each group separately from the fixed effects.

[5]For discussions of the error term assumption and its implications on substitution patterns, see Hausman and Wise (1978), Berry, Levinsohn, and Pakes (1995), Nevo (2001), Petrin (2002), Bajari and Benkard (2003), and Berry, Levinsohn, and Pakes (2004).

in the list, households make a decision on the purchase quantity for each good as a function of expected prices and other state variables.

We model the choice of goods purchases in the first stage using a continuous-time duration model that accounts for heterogeneity across households in goods purchasing behavior. The composition of basket goods in a shopping list is then characterized by the incidence of each good's being included in the shopping list. Conditional on these estimated incidences, the quantity of each good that a customer plans to buy is predicted using a regression of realized quantities on actual prices, brand fixed effects, and other relevant variables. In specifying both of the basket components (goods and quantity), the choices are assumed to be made independently among different goods.[6]

For the unobserved price expectation, the model posits that households develop price knowledge based on their experiences from past shopping visits. Our modeling assumption is that the expected price is a weighted average of the past prices that a household has observed or paid during previous store trips. Price expectation defined in this way depends on good purchases and shopping patterns of each household. The richness of household level information on shopping trips and purchases allows us to construct household-, time-, store-, and good-specific price expectation.

Given *ex ante* shopping lists and price expectation, households are assumed to make a store choice following the store choice model described in the previous section. In the store choice, each household chooses the store that maximizes the utility from shopping for the expected shopping basket. In the following sections we provide further details on the three components that determine expected basket costs: price expectation, goods purchase incidence, and quantity of purchased goods.

---

[6]This assumption is not problematic for the goods that are irrelevant to each other such as milk and laundry detergent. However, there may exist a set of goods for which their purchase and consumption are interdependent, such as hotdog buns and ketchup goods. In our data there is only one case of such dependency in purchase decisions (toothbrush and toothpaste), and we ignore this issue in this study.

## 3.1 Price Expectation

In our price expectation formation model, we make several behavioral assumptions on consumers. First, households form price expectation based on past prices from their recent trips to stores, which depends on each household's shopping behavior and, therefore, accounts for heterogeneity in price knowledge. Second, given that grocery stores carry a vast array of products and consumers do not remember most of the prices given time constraints for each shopping trip, our simplifying assumption is that the past prices of products are only kept in memory when consumers have purchased the goods categories into which those products fall.[7] Lastly, memories of past prices fade over time, and thus, price expectation for each product is closest to the prices from the most recent trips to the store. This approach of constructing price expectation is similar to the reference price model, first introduced by Winer (1986), in that consumers adaptively formulate price beliefs or forecasting rules for price based on personal histories of purchases or information process. More recently, Briesch, Chintagunta, and Fox (2009) pointed out that using smoothed past prices, which is in line with our measure of price expectation, can be an alternative to assuming rational expectations of consumers.

Based on these behavioral assumptions on the formation of price expectation, expected price is defined as the weighted average of past prices at the store where the weights exponentially decline with remoteness in time from the current period. Formally, the expected price of product $k$ of good $c$ for household $i$ and store $j$ at time $t$ is written as

$$\bar{p}_{ijkt} = \sum_{\tau < t} w(\tau; t, \kappa) d^s_{i\tau}(j) d^c_{i\tau}(c) p_{jk\tau}, \tag{4}$$

where $p_{jk\tau}$ is the actual price of product $k$ at store $j$ at time $\tau$, $d^s_{i\tau}(j)$ is the indicator of store visit which equals 1 if the household visits store $j$ and 0 otherwise, and $d^c_{i\tau}(c)$ is the good purchase indicator which takes 1 if the household buys any product of good category $c$ at time $\tau$ and 0

---

[7] If shoppers recollect the prices of all products regardless of whether they were shopping for the good during a store visit, price expectation becomes almost homogeneous among households since most households visit stores as frequently as once per week. For example, if two customers visit the same store almost every week, their price beliefs for each good will be almost identical even though they have purchased different goods for each trip. In this case, heterogeneity in price expectation stems only from the frequency of store visits no matter which goods they have purchased in past shopping trips. Our general model setting allows this simple case too.

otherwise. The weight function $w(\tau; t, \kappa)$ is specified as

$$w(\tau; t, \kappa) = \frac{\exp(-\kappa(t - \tau))}{\sum_{\tau' < t} exp(-\kappa(t - \tau'))},$$ (5)

where the parameter $\kappa$ may be interpreted as the degree of memory decay. A high (positive) value of $\kappa$ puts greater weight on more recent prices, and the coefficient of zero equally weights the past prices (thus, price expectation becomes a simple average of the past prices). Price expectation specified in this manner not only captures the short-term variation of store prices by placing larger weights on more recent prices, it also accounts for heterogeneity in shopping behavior since it is based on the goods purchases and the store choices of each household in the past.[8]

This approach to deal with unobserved price beliefs can be viewed as a generalization (or in the least a robust-check) of the price expectation used by Bell, Ho, and Tang (1998) and Briesch, Chintagunta, and Fox (2009). They assume that consumers have common knowledge on the price distribution for each store up to the first and second moments. Then, they define the expected price of each product as the average of store prices over the data period. Their specification of price expectation implies that shoppers visit each store frequently enough and acquire price information from various sources to the extent that the price process is known up to the first moment of the price distribution. In such a way, their price expectation reflects a long-term variation in prices including future prices as well as past prices. In our setting of price expectation, given in equation (4), their expected price would be written as

$$\bar{p}_{ijkt} = \sum_{\tau < \infty} w(\tau; t, \kappa = 0) p_{jk\tau}.$$

That is, $d_{i\tau}^s(j) = d_{i\tau}^c(c) = 1$ for any store $j$ and good $c$ where the memory decay coefficient in the weight function is set at zero. Note that in this case price expectation is identical among all households and thus leaves out heterogeneity in price beliefs.[9]

---

[8]Since past prices are embedded in price expectation only when the good is bought, customers have no price knowledge if they have never purchased the good at the store (that is, $d_{i\tau}^c(c) = 0$ for any $\tau < t$). In this case, we set the initial price expectation based on past shopping trips assuming they observed the prices although they did not purchase the goods.

[9]In an alternative setting we loosen their assumptions and estimate the store choice with the price expectation

## 3.2 Goods Purchase Incidence

We model the goods purchase decision such that it consists of a set of sequential purchase incidences in which consumers visit stores to buy some products of a specific good at randomly selected times. We characterize the composition of basket goods in terms of the likelihood of good purchase incidence (so called "hazard" in the duration model literature) as a function of the length of the elapsed time since the most recent good purchase. The hazard rate is also determined by other state variables that directly or indirectly affect purchase decisions at a given point of time. We briefly introduce the duration model of goods purchases, adopting the standard notation in the duration model literature (e.g., Cox (1972), Heckman and Singer (1986), Lancaster (1990), and Kalbfleisch and Prentice (2002), among others).

Let $T$ be a random variable of the duration of a state, or the elapsed time since the most recent purchase. Conditional hazard $h(\tau)$ is the instantaneous probability of leaving the no-purchase state (i.e., the occurrence of a new purchase incidence) after time $\tau$. Given $x_{it}$, a vector of observed state variables for household $i$ at time $t$, the hazard rate can be written as

$$h_i(\tau|x_{it},\mu) = \lim_{\triangle \to 0} \frac{Pr\left[\tau \leq T < \tau + \triangle | T \geq \tau, x_{it}, \mu\right]}{\triangle}, \tag{6}$$

where $\mu$ is a vector of parameters. Using the specification proposed by Cox (1972), the conditional hazard (6) is given as

$$h_i(\tau|x_{it},\mu) = h_{i0}(\tau)\phi(x_{it},\mu). \tag{7}$$

In this specification the conditional hazard is assumed to be proportional to the two components on the right-hand side. The first component $h_{i0}(\tau)$ is the baseline hazard as a function of no-purchase spell $\tau$ only, and the second component $\phi(x_{it},\mu)$ is a function of observed state variables $x_{it}$ at time $t$. The baseline hazard $h_{i0}(\tau)$ accounts for heterogeneity in shopping behavior across households (particularly the frequency of goods purchases). In equation (7), the parameters in the second component of the hazard is estimated by specifying the functional form of $\phi(\cdot)$ whereas

---

only under the assumption that $d_{i\tau}^c(c) = 1$ for any $c$. That is, households memorize prices when they visit the store but regardless of whether they purchase the goods or not. We find the estimated parameters for expected basket costs with this price expectation in the store choice are actually positive, which is not reasonable. The estimation results for this case will be provided upon request.

the household-specific baseline hazard is not parameterized. Specifically, for our estimation, we use the specification $\phi(x_{it}, \mu) = \exp(x'_{it}\mu)$.

In the duration model, the time-variant state variables $x_{it}$ for each grocery good may include information on consumption and inventory levels, purchases of other goods, and individual households' shopping patterns and demographics.[10] Since the consumption and inventory levels are not directly observed in the data, we construct these variables based on the observed purchases (quantity and frequency) of the good for each household considering the specific features of each good such as average shelf lives and approximate storage costs.[11] More details on how we construct the consumption and inventory variables are provided in Appendix B.

Including the variables of consumption level and purchases of other goods in $x_{it}$ accommodates exogenous shocks to the need of purchasing the good of interest. For example, a household that recently consumes or purchases a large amount of other goods may be more likely to buy the good of interest soon after previous purchase. The variables for shopping behavior are constructed based on the observed grocery shopping patterns of each household, and these include the average frequency and dollar spending of recent shopping visits, preference for weekend shopping, the average time interval between goods purchases, and seasonal dummies. Here our underlying modeling assumption is that such shopping behavior is rather persistent for each household and does not vary over the sample span, so the approximation of such shopping pattern is possible using observed data.

In our model setting, goods purchase decisions themselves are not affected by the expected prices of the products (while the quantity choices depend on the expected prices given the purchase decisions). We assume that the effect of price changes on the need or consumption of each good is negligible or none. For example, people do not increase their average consumption level of carbonated beverage or laundry detergent or purchase them more frequently when they

---

[10]The variables to control for shopping patterns are also based on overall grocery shopping behavior, not just related to the good of interest.

[11]It may be natural to consider the consumption level as a decision variable in store choice as in Hendel and Nevo (2006), who focused on stockpiling behavior in one particular grocery good. However, since store choices involve a wide variety of goods and the consumption decision is an object that should be understood in a dynamic framework, it is not straightforward to allow consumption to be endogenously determined. For this reason, in our setting, we assume an exogenous consumption level given the endogenously determined purchase incidence decision.

expect low prices. In addition, the effect of expected prices on goods purchases will be only tenuously identified in our setting because price expectation is updated after households purchase the goods and thus the expected prices do not change during no-purchase periods.[12]

Specifying a shopping list at the goods level instead of at the product level makes the model flexible enough to allow substitution between different products of the same good to take place inside stores. If households have strong loyalty to specific brands or sizes of each good, they will always plan to buy specific products, and thus a product-level shopping list would be more appropriate. However, it is well documented in both the economics and the marketing literature (e.g. Kumar and Leone (1988), Blattberg, Briesch, and Fox (1995), Pauwels, Hanssens, and Siddarth (2002), and Hendel and Nevo (2006)) that the consumer choice of brand or size is determined to a large extent by in-store promotions. Therefore, it seems reasonable to define a shopping list at the goods level.[13] In contrast, Bell, Ho, and Tang (1998) specify a planned shopping list at the product level and the list consists of the products that were actually purchased at the store. This approach may be more reasonable when households actually have strong loyalty to specific brands or sizes of each good, but it leaves out the possibility that the purchased products were not contained in the planned shopping list and did not influence store choice decisions.[14]

## 3.3  Quantity of Purchased Goods

The third component that determines expected bundle cost is the quantities of the goods that shoppers plan to purchase. Households choose the purchase quantity for each basket good based

---

[12]If price expectation reflects price information obtained from the sources other than goods purchases such as price advertising, expected prices will evolve during no-purchase periods, and this would affect the timing of the upcoming goods purchases. However, we presume that the role of price advertising is limited in the context of store choice decisions. Bodapati and Srinivasan (2006) document that only a small fraction of customers are influenced by price advertising in their store choices. Our data also show that advertising is placed for less than 10 percent of the products sold by each store in any given week. Bell, Ho, and Tang (1998) also did not use price advertising information.

[13]Modeling a shopping list at the goods level is also supported by Block and Morwitz (1999) who found that consumers write either grocery categories or specific products on their shopping lists, but 77% of the listed items are at the category level only.

[14]There have been several studies documenting that realized purchases could differ from planned shopping lists using survey or interview data. For example, Kollat and Willett (1967), Block and Morwitz (1999), and Bell, Corsten, and Knox (2011) document that more than 50% of items (or 20% of categories) purchased were not planned.

on their price expectation and other relevant states such as shopping frequency, consumption and inventory level for the good. Since neither the price expectation nor the planned purchase quantity is observed in the data, we have to infer how much consumers plan to buy based on the observed quantity choices at stores. Then, to predict the expected quantity, the shelf prices that households observe inside a store are replaced by their price expectation in the equation of the quantity choice. Specifically, only the coefficients in the quantity equation are estimated by regressing the realized quantity on actual price taking into account brand and household fixed effects, seasonal dummies, and demographic variables. The estimation equation for the quantity purchased is

$$q_{ijkt} = W_{ikt}\beta^q - \alpha^q p_{jkt} + Z_{it}\gamma^q + \delta_i^q + \delta_b^q + \delta_t^q + \epsilon_{ijkt}, \tag{8}$$

where $q_{ijkt}$ is the quantity in volume of product $k$ purchased by household $i$ at store $j$, $p_{jkt}$ is the unit price of the product, $Z_{it}$ is a vector of demographic characteristics, and $W_{ikt}$ is a vector of state variables such as the shopping frequency of the year and the consumption and inventory levels for the good of interest.[15] We also include the dummies for household, brand, and times (year and month), denoted by $\delta_i^q$, $\delta_b^q$, and $\delta_t^q$, respectively. Again, when we predict the planned quantity based on the estimated parameters, we use the expected prices as prices. Here our modeling assumption is that households make the planned quantity choices as if the expected prices were the actual prices.

To deal with the potential issue of endogeneity in prices, actual prices in the regression are instrumented by the average prices of each product in the same supermarket chain stores in nearby cities. As in Nevo (2001), if stores operated by the same supermarket chain have a similar cost structure (thus correlated with the price in equation (8)) and if the city-specific valuation of the price is independent across different cities, the prices at the same chain stores in other cities can be valid instrument variables.[16]

---

[15]The purchase quantity of each product has a common measure within the same good and the price is defined for each measure unit. For example, the measure of quantity for milk good is gallon and the price is defined for each gallon.

[16]As pointed out by Erdem, Imai, and Keane (2003), endogeneity in prices of frequently purchased consumer goods may be more attributable to the omitted variables such as consumer inventories than to aggregate demand shocks. Since we control for inventory level in our quantity estimation, we use the average prices in nearby cities for dealing with the endogeneity problem that may arise from unobserved demand shocks.

Bell, Ho, and Tang (1998) also acknowledge that the realized quantity may differ from the expected quantity in a shopping list. But they argue that there is only a slight improvement in their estimation results when considering the discrepancy between realized and planned quantity. There can be two reasons for the limited role of quantity choice in their study. First, if the estimation of quantity choice poorly fits the data, the prediction of planned quantity is not accurate. Second, this may be because expected prices in their model do not allow for heterogeneity in price knowledge as discussed above.

## 3.4 Expected Bundle Cost

Given the three components of a shopping list as described above, the expected basket cost is defined as the sum of the expected spending on each basket good. For each basket good, the expected spending is the average of the expected costs of the products of the good. Then, the expected bundle cost is the weighted average of the expected spending on each good with the weight of the good purchase likelihood (i.e. hazard rate). Formally, the expected basket cost is written as

$$Exp_{ijbt} = \sum_{c \in \mathcal{C}(b)} h_i^c(\tau|x_{it}, \mu) \sum_{k \in \mathcal{G}_{ic}} \frac{\bar{p}_{ijkt} E(q_{ijkt}|\bar{p}_{ijkt}, W_{ikt}, Z_{it})}{n(\mathcal{G}_{ic})}, \tag{9}$$

where $\mathcal{C}(b)$ is the set of goods contained in shopping basket $b$, $h_i^c(\tau|x_{it}, \mu)$ is the hazard rate for good $c$ after a duration of time $\tau$ since the most recent good purchase, and $n(\mathcal{G}_{ic})$ is the number of products included in the consideration set $\mathcal{G}_{ic}$. $E(q_{ijkt}|\bar{p}_{ijkt}, W_{ikt}, Z_{it})$ is the expected quantity for product $k$ at store $j$ as a function of expected price $\bar{p}_{ijkt}$ and other variables.

In calculating the expected spending for each good, we assume that households only consider the products that they can possibly buy at stores. Based on this assumption, the expected cost only accounts for the products in the consideration set of each household for the good. The consideration set $\mathcal{G}_{ic}$ includes the products that the household ever purchased in the sample period. For example, if a household has only bought Coke, Pepsi, and Dr Pepper, the expected spending of soft-drink good for this household is defined by the expected costs of these three brands. Since the household level data cover five years, the consideration set defined in this way comprehensively captures the substitute brands and sizes for each household.

Restricting the set of products to a consideration set for each household provides another source of household heterogeneity in price expectation. The variation in price expectation arising from this is substantial because the grocery stores carry a vast array of products for each good category and their pricing strategies may differ across these products among competing stores. For example, if the prices of Pepsi and Mountain Dew products are relatively low and those of Coke and Dr Pepper are high in a store compared to other stores, depending on which products households consider buying, the expected cost of soft-drink good for each store can vary across both households and stores.

We note that the expected basket cost, $Exp_{ijbt}$, potentially has substantial variation across households and stores, and over time. The sources of such variations are multifold. First, the price expectation for each product builds on the household's shopping history. Second, the consideration set for each good differs across households depending on their choices of brands and sizes. Third, heterogeneity in the decisions of the expected quantity also creates variations in the expected spending for each good across households. Last, the good purchase incidence predicted from the shopping behavior of each household is another source of variation in the expected basket cost.

## 4    Data

In our application we use scanner data collected by IRI.[17] The data set is in two parts. The first data set is store level data containing weekly store sales and the second is household level data containing the weekly purchases and store visits of individual households. The data set covers five years from January 2003 to December 2007 and includes 30 grocery goods, which consist of 17 food and beverage goods (e.g., carbonated drinks, coffee, cereal, frozen meals, peanut butter, soup) and 13 non-food household goods (e.g., diapers, facial tissues, laundry detergent, shampoo, razors, toothpastes). The data used in our analysis are drawn from seven stores in a small city in Massachusetts and the sample stores belong to four different supermarket chains. The store level data include weekly price and quantity for each product sold in each store at the

---

[17]Further details on the data sets are provided in Bronnenberg, Kruger, and Mela (2008). We thank IRI for making this data set available to us. All analysis using the data in this paper is by the authors and not by IRI.

UPC (universal product code) level, covering a much larger number of stores.[18] The household level data contain information on individual purchases of each household at each store. We restrict the sample to about 1,700 customers who report their purchases persistently enough to satisfy the minimal reporting requirements set by IRI.[19] We use the last two years of the data, 2006 and 2007, for the estimation of store choice, while the first three years are also included in constructing price expectation and estimating the choice of goods and quantities.

Since the price variable in the household level data is the weekly average of store prices, it may not be the same as the price that was actually paid by the customer. This may generate measurement errors in prices if customers redeemed retailer coupons or if the store prices change within a week. However, the cases in which coupons are offered by stores are only about 0.05 percent in the store level data. Therefore, it would be mostly the case that a discrepancy between the prices in the data and the actually paid price may occur if the shelf prices vary within the week.

We have information on demographic characteristics for each household, such as income, age, home ownership, dummy variables for single male and single female, number of children, and an indicator for full-time working female. We also have information on the location of each household and the sample stores with latitude and longitude.[20] We computed the shortest driving distance between the residence of each household and the sample stores using the Google Maps API with the location information. Using the detailed location information is unique to our data set. Most previous studies on store choice approximate the trip distance from the zip code or census block information. Given it is well documented that the spatial distribution of stores is one of the key determinants in store choice (see e.g. Smith (2004), Thomadsen (2005), Davis (2006), Ellickson and Misra (2008), Briesch, Chintagunta, and Fox (2009), Houde (2012), and Orhun (2013)), using an accurate measure of trip distance is crucial in estimating store choice consistently.

---

[18]The original store level data cover about 2,500 stores, including the seven stores used in our study, in 50 U.S. cities.

[19]The persistency of reporting is evaluated by IRI every year. We only include the households who meet the reporting requirements continuously for the full participating period. For example, if a customer meets the criteria in 2005 and 2007 but does not in 2006, this customer is dropped from the sample.

[20]For confidentiality reasons, the exact location of each household is disguised by a trifling error (about 0.1 mile).

The geographic area where the sample customers and stores are located is a small urban area. This area is appropriate for the study of store choice because the consumers in this market do not have non-traditional grocery stores nearby, such as supercenters or warehouse clubs which are not included in the data. The non-traditional stores provide not only grocery products but also a large variety of general merchandise goods. The presence of such stores in a market can make a store choice study complicated because the purpose of grocery shopping can be confounded by purchases of non-grocery products. The supermarket chains included in the data are all traditional supermarkets, and the closest supercenter supermarket (Walmart or Target) is 19.6 miles from the sample customers on average, whereas the average distance between the sample customers and the stores in the data ranges between 2.9 and 5.6 miles.

Using the consumer level data containing a subset of households in this area may raise a potential concern about sample selection bias. Households who participate in the data collection may not be representative of the market since they need to commit a certain amount of time and effort in reporting their grocery purchases. For example, a household whose opportunity cost of time is relatively high may be naturally excluded from the sample customers. Table 1 presents the comparison of demographic statistics between the sample households in our data and the population in this area surveyed by the U.S. Census Bureau for a subset of the data periods, between 2005 and 2007. The sample households are about similar to the population in the distribution of household income but the sample distribution has slightly thinner tails than the population. However, the median income is about the same at $43,000 between both statistics, which is not reported in the table. The household sample in our data has a slight bias in selection towards the middle-aged or elderly customers (perhaps, the reason why the data contain a smaller number of full time working women). The proportion of single female is similar between our data and the population. However, single male households are far less selected than the population in this area.

Table 2 provides descriptive statistics of the household level data. The frequency of store visits tells that the average consumers are making shopping trips almost once per week.[21] But

---

[21]The frequency of store visits might be underestimated if households visit a store multiple times in a week because these are not distinguished in the data. However, the supplementary data set, provided by IRI, that contains all the store visits – not confined to the purchases of the 30 goods – show that about 70 percent of visits

there is a large variation in the trip frequency across households. For example, shopping frequency is negatively correlated with dollar spending per trip (correlation is -0.5), implying that shoppers with a large shopping basket tend to visit stores less often than those with a small basket. Consumers visit about four different stores per year on average. But the Herfindahl–Hirschman Index (HHI) of store choice, which measures the extent to which the store choice of each household is concentrated among the stores, amounts to 0.5, based either on number of visits or dollar spending. This suggests that, on average, store choices are concentrated on nearly two different stores (see similar empirical findings in e.g. Gauri, Sudhir, and Talukdar (2008) and Briesch, Dillon, and Fox (2013)). Households buy three different goods on average, among the 30 goods categories, for each store visit.

Figure 1 depicts the distribution of the store HHI of households, computed based on the number of store visits and dollar spending. It shows that the concentrations of store choices substantially differ among households. It also shows that, for most households, neither spending nor store visit is concentrated on one or two stores only.[22] Examining the extent to which various factors in store choice affect switching behavior is left to an empirical question.

In our application, we assume that transportation cost is a linear function of the driving distance between the residence of households and the stores. Figure 2 shows the choice of stores by the rank in distance among alternative stores. The nearest store is the most popular choice, both in the number of store visits and dollar spending, which reinforces the conventional wisdom that transportation cost is crucial in store choice. However, given that the difference in distance between stores in two consecutive ranks is only half a mile on average, it is not straightfoward to what extent the distance affects store choice decisions.

Most households do not visit all the sample stores during the sample period. The reason why they never visit certain stores is not obvious. It may be because the utility from visiting those stores does not exceed those from visiting the stores they usually choose or because they do not know about the stores and thus those stores are simply not in their choice set. Including

are made only once to each store in a week. Considering that the trip data cover all the store trips made for any goods, the case of multiple visits would be much less problematic for 30 goods we use for our estimation.

[22]Alternatively, the HHI can be calculated based on the choices of supermarket chains instead of stores. The distribution of the chain HHI barely changes compared with that based on the store choices.

these never-visited stores in the choice set therefore may lead to inconsistent estimates of demand parameters. For example, suppose a household has little knowledge about a store for some reason and never visited the store, but the model includes this store in the household's consideration set. If the prices at this store are low enough that the utility from visiting this store would be higher than those from other stores, the estimated price elasticities will be biased toward zero. For this reason, we exclude these never-visited stores from the choice set of each household.[23]

Table 3 presents the summary statistics of the store level data. According to the non-disclosure agreement with IRI, the chain names of the sample stores are disguised. The first two columns show the percentage frequency of price discounts that are larger than 5 percent of regular prices estimated by IRI with a proprietary algorithm, and the number of different products carried by each store, respectively. There is a substantial variation in price promotions across stores and the variation is much larger for some goods than others. It is noteworthy that price promotions and assortment sizes differ across stores within the same chain. The average driving distance between the sample customers and each store ranges from 3 miles to 5.6 miles.

# 5  Estimation and Results

## 5.1  Estimation of Expected Basket Costs

In the first stage, price expectation of each household is constructed based on price knowledge from previous shopping trips as described in Section 3.1. Since our model allows the expectation of store prices to be developed or revised when the household buys certain good category at the store, the expected prices can be missing prior to the first-ever purchase of the category at the store. To deal with this censoring problem, we use the first three years of our sample only for the purpose of constructing price expectation and exclude them in the store choice estimation. The parameter of memory decay rate, $\kappa$ in the weight function of price expectation (5) is estimated based on the full model, using a profiling estimation method. That is, we fix the parameter $\kappa$ at a specific value and estimate the model of the store choice. By repeating this procedure with

---

[23]If a store loyalty variable is included as an alternative specification such that it accounts for habitual choice behavior based on the observed store choices in a certain length of an initial period of data (e.g., Bell, Ho, and Tang (1998)), this variable would absorb this biased selection in store choices.

different values of $\kappa$ we obtain the profiled estimates of other parameters. We then pick the value of $\kappa$ that maximizes our objective function of estimation and evaluate other parameters at this estimated value of $\kappa$ (see further details in Section 5.2). The estimated value of $\kappa$ is 0.032, which indicates the memory decay is relatively small but it is statistically different from zero. The weight function at $\kappa = 0.032$ is only moderately steep and it decreases gradually. With this weight function, past prices of the recent three months are given about 40 percent of the weights out of the whole one year period. This implies that past prices embedded in price expectation are smoothed over a somewhat long period rather than only recent prices being picked out.

Since price expectation in our model reflects store trips and goods purchases at different times depending on each household's shopping pattern and history, expected prices account for household heterogeneity in price beliefs. Figure 3 illustrates the price expectation of four randomly picked households for a popular cereal product over one year. The solid line shows that the store prices of this product stay at regular prices for most of the time and drop sporadically with temporary price cut for a short duration of 1 or 2 weeks. Most of the grocery items in our data follow similar price patterns. Price expectation of each household changes in a different manner depending on its shopping pattern. If a household shopped at the store while price discounts were offered, for example, the expected price declines accounting for the update in price information. Also, the expected price tends to rise as the regular price increases or price promotions are offered less frequently for the second half of the year in Figure 3.

Next, we estimate the likelihood of a good purchase incidence in a shopping plan separately for each good using the hazard model as described in Section 3.2. We estimate the hazard rate (i.e. likelihood of purchasing incidence) as a function of log of inventory and consumption levels, recent shopping patterns, and demographic information. The variables we use to control for recent shopping patterns are trip frequency, spending per trip, weekend-shopping preference, holiday dummies, and time intervals (in week) between recent goods purchases. Demographic variables include log of household income, family size, and indicators of marriage, pet ownership, and house ownership. In the estimation we also add the recently purchased quantities of other goods to control for unobserved exogenous shocks to the good purchase of interest.

Table 4 reports the estimation results for six selected goods. The coefficients represent the

impact of each variable on the good purchase incidence. Most of the variables are statistically significant and the signs are intuitive. The negative coefficient of inventory level implies that, the larger amount of the good in a pantry (or the lower hazard rate), the lower probability of purchasing the good, that is, households wait longer until the next purchase. Higher consumption is associated with a shorter purchase interval. Shoppers with a high visit frequency and large dollar spending for each trip tend to have a high chance of the upcoming purchase. Purchase incidence seems closely related to holiday seasons. Purchases of soft-drink tend to be more likely during the holidays, while less so for such goods as cereal, laundry detergent, and paper towels.

Figure 4 shows the survivor functions for different levels of household inventory for the six categories of goods, holding other covariates fixed at their mean values. The survivor function $S(\tau|\bar{x})$ is the probability that the duration until the next purchase exceeds $\tau$ given $\bar{x}$ where $\bar{x}$ denotes a vector of sample means of observed state variables such as consumption and inventory levels, purchases of other goods, and individual households' shopping patterns and demographics. Specifically, we compute the survivor functions for the 10th and 90th percentiles of inventory levels of the sample households. The survivor functions are monotonically decreasing with the elapsed time since the most recent purchase. The declining speed of survivor functions varies among different goods since the purchase frequencies are different. For example, a subsequent purchase of blades can occur after more than a year with a positive probability, whereas the survivor functions for milk reduce to zero within about five to six weeks. More importantly, a higher level of household inventory at the moment of the recent purchase is associated with a higher survival probability at any point in time. This implies that the households wait longer until the next purchase with a higher level of inventory level.[24]

Next, we estimate parameters of the quantity choice equations as described in Section 3.3 to predict the planned quantity in the shopping list estimated above. Table 5 reports the estimation of quantity choice for the six goods. Since most goods are sold in numerous sizes of package, the price variable in the estimation is rescaled to a standard size for each good. The price coefficients are negative for all goods, and the magnitudes of the coefficients in absolute

---

[24]Survivor functions for other goods are consistent with these findings. The figures of the survivor functions for other goods can be provided upon request.

terms increase, compared to the OLS estimates (which we do not report here), when the price is instrumented. This suggests the endogeneity of the observed prices. Allowing for brand and household fixed effects substantially improves the fit of the quantity estimation. Using these coefficient estimates, we predict the planned quantity as described in 3.3 as an input to construct the expected basket cost.

Finally, combining estimates on the expected prices, goods purchase incidences, and the planned quantity choices, we construct the expected basket cost using the formula (9) as described in Section 3.4. Table 6 reports the expectation errors in the expected basket cost across different households grouped by their trip frequencies. Expectation error is defined as the root mean square of the difference between the predicted expected cost and the observed actual spending for the basket goods. The actual basket spending is based on actual store prices, which households do not observe *a priori*. Expected basket cost in the full model is based on all the sources of heterogeneity across households as described above, whereas the non-heterogeneity model does not allow for household-specific shopping histories and consideration sets. The full model in the table shows that frequent shoppers have relatively smaller expectation errors than infrequent shoppers. The frequency is defined as the number of shopping trips to the store during each year. It can be inferred that price expectation becomes more up to date with frequent past store trips. The expectation errors based on the non-heterogeneity model are larger than those in the full model, on average, for all household groups. More importantly, the variation in expectation errors across different household groups does not show the same pattern with the full model.[25]

## 5.2 Store Choice Estimation

Using the estimated expected basket cost in the first stage as a regressor, we estimate the store choice model we develop in Section 2. We use a simulated MLE method based on the multinomial logit discrete choice model with the store choice probabilities given by (3). For any candidate

---

[25]Note that the expectation error in the full model for the third group is somewhat smaller than that for the fourth group (i.e., they are not exactly proportional to the shopping frequency). The expectation errors presented in Table 6 may not fully capture household heterogeneity in price beliefs because the trip frequency is merely one of the various shopping patterns that characterize each household.

value of $(\theta, \kappa, \nu_i)$, where $\theta = (\alpha_0, \beta_0, \{\gamma_{0j}\}_{j \in J}, \Lambda, \Sigma)$ is the vector of the parameters to estimate, we write the conditional (on $\nu$) choice probability equation from the logit model as

$$s_{ijb_{it}t}(\theta, \kappa, \nu_i) = \frac{\exp(\beta_i Dist_{ij} - \alpha_i Exp_{ijb_{it}t}(\kappa) + \gamma_{ij})}{1 + \sum_{j' \in J} \exp(\beta_i Dist_{ij'} - \alpha_i Exp_{ij'b_{it}t}(\kappa) + \gamma_{ij'})},$$

where $b_{it}$ denotes the planned basket chosen by household $i$ at time $t$, $(\alpha_i, \beta_i, \gamma_{ij})$ are the heterogeneous coefficients defined in the equation (2), and $Exp_{ijb_{it}t}(\kappa)$ denotes the predicted basket cost expectation evaluated at the memory decay parameter $\kappa$. For the estimation, we use the simulated maximum likelihood assuming the random coefficients are normally distributed where $\nu_i$ is independent normal errors (i.e. $\Sigma$ is a diagonal matrix of variances).[26] The (profiled) maximum likelihood estimator is then given by

$$\hat{\theta}(\kappa) = \text{argmax}_\theta \ \log L(\theta, \kappa) \equiv \log \prod_{i=1}^N \prod_{t=1}^T \left( \frac{1}{R} \sum_{r=1}^R \left( \prod_{j=1}^J \left( s_{ijb_{it}t}(\theta, \kappa, \nu_i^r) \right)^{y_{ib_{it}t}(j)} \right) \right), \qquad (10)$$

where $\nu_i^r$ is the $r$-th simulation draw from the normal distribution, $R$ is the number of simulation draws for $\nu_i$, and $y_{ibt}(j)$ is the indicator function such that

$$y_{ibt}(j) = \begin{cases} 1 & \text{if household } i \text{ at time } t \text{ chose the store } j \\ 0 & \text{otherwise.} \end{cases}$$

We use a crude frequency simulator with 100 pseudo-random draws of $\nu$ for each household. Using the profiled likelihood in (10), we then estimate the memory decay parameter $\kappa$ such that

$$\hat{\kappa} = \text{argmax}_\kappa \ \log L(\hat{\theta}(\kappa), \kappa) \qquad (11)$$

and finally obtain the estimator for $\theta$ as $\hat{\theta}(\hat{\kappa})$ where $\hat{\kappa}$ is obtained from (11).

Using this ML method, we estimate the store choice model for different cases of expected bundle costs. First, the full model employs the expected basket cost reflecting all sources of household heterogeneity, as described in Section 3.4. Then, we estimate the cases in which each

---

[26]See e.g. Train (2009) for the simulated ML methods.

of the heterogeneity sources is eliminated in the expected bundle cost of the full model. This allows us to examine to what extent each of the components incorporated in the expected bundle cost makes a difference in the estimation of store choice.

The estimated parameters of the full model are reported in Table 7. The mean coefficients of marginal utilities are presented in the first column, and the next columns present the estimates of household heterogeneity around these mean values. The mean coefficients are all statistically significant and take the expected signs. The base group for the interactions with income and age is set as households older than 65 with annual income less than 15,000 dollars. Households with higher income compared to the lowest income group are less price-sensitive, but not proportionally, and older households are more sensitive to the price level. The marginal disutility of travel distance increases with income level, which is intuitive given that hourly wage is expected to be positively correlated with income level. Young households tend to be less unwilling to travel farther for shopping compared to those older than 65, although customers aged less than 34 prefer visiting nearby stores. Estimates of the standard deviations are all statistically significant and their magnitudes are relatively large compared to mean coefficients. This suggests that there exists non-negligible amount of individual heterogeneity beyond what is explained by the observed demographic characteristics.

Table 8 presents the estimates from using different measures of the expected bundle cost. The first three columns show the results based on our *choice-based* price expectation model. The second and third columns are the cases for which the estimation of either purchase quantity or goods purchase incidence is removed, respectively, in constructing the expected bundle cost. The next three columns are based on price expectation that is modeled as not (or only partially) allowing for heterogeneity in price beliefs.[27] The price expectation used in the fourth column is a simple average of store prices over the past one year, regardless of the shopping history of individual households, but the expected bundle cost is estimated reflecting the difference among households in the product consideration set. The fifth column is using the measure of the expected bundle cost that ignores both household-specific shopping history in price expectation

---

[27]In the three cases of homogenous price expectation models, the measures of the expected bundle costs are still calculated using the estimated goods purchase incidence and purchase quantity. Removing these two components of the basket model yields a similar comparison as in the choice-based price expectation.

and heterogeneity in product consideration set.[28] In this case, there is no component accounting for heterogeneity in price expectation and the expected prices are all the same across households. The last column is similar to the fourth column in constructing price expectation, but the shopping list in the basket model is defined at the product level as modeled by Bell, Ho, and Tang (1998).[29] Therefore, comparing this case with the fourth column may give a hint on the role of modeling the goods-level shopping basket.

The comparison of the first three columns suggests that the estimation of goods purchases incidence and purchase quantity in the first stage makes only a slight difference. However, when the purchase quantity estimation is removed, the heterogeneity around the mean price coefficient among different income groups almost dissipates and becomes statistically insignificant compared to the other two cases of the choice-based price expectation. The fourth and fifth columns show that ignoring heterogeneity in price expectation substantially changes the magnitudes and signs of the coefficients on the expected bundle cost. When price expectation does not properly account for household heterogeneity, the mean price coefficients are biased toward zero and become even positive. When a shopping list is defined at the product level, such biases in the price coefficient seem to be reduced compared to the fourth column. However, when compared with the first three columns, the variation of price sensitivity by age shows a different pattern. In this case, younger customers are more price sensitive than aged shoppers whereas it is the other way around in the cases of the choice-based price expectation. The negative sign of the mean distance coefficient in all cases supports the finding from previous work that consumers have a substantial amount of disutility from traveling to a distant store.

Since the coefficients estimates for the basket cost given in Table 8 only indicate the marginal changes in terms of the latent utility, we translate them into price elasticities. Table 9 presents the own elasticities of demand with respect to the expected basket cost for each store. The elasticities are calculated as the percentage changes in the predicted market share of store choices when the expected basket costs of households increase by one percent. The estimates based on

---

[28]It is not possible by construction to remove heterogeneity in the consideration set alone because the past prices of never-bought products are not defined if households consider prices of all products instead of only those in their consideration set.

[29]In our replication of the work of Bell, Ho, and Tang (1998), we define the expected price as the average store price over one year whereas theirs is the average of the entire data period (two years)

the choice-based price expectation yield negative own price elasticities and the magnitudes are quite similar among the three cases. This implies that using models for the goods purchase incidence and the purchase quantity of planned goods in the expected bundle cost does not matter much in the store choice estimation, and therefore a parsimonious approach of using the realized purchases as a proxy for planned shopping list may be used. On the other hand the fourth and fifth columns suggest that, when heterogeneity in household price expectation is ignored, the estimates of own price elasticities can be severely biased toward zero and even take the opposite sign. The comparison between the fourth and the last columns shows that defining a shopping list at the product level somewhat mitigates these biases in own price elasticities but the magnitudes of the elasticities are fairly small compared to those obtained by our approach.

# 6    Conclusion and Discussion

This paper proposes a household-level store choice model that reduces the potential biases arising from mis-measuring unobserved price expectations and shopping lists of households. We specify the household-level price expectation based on shopping history and the set of products chosen by each household for each good. This captures the effects of the short-term variations of store prices on each household's price knowledge and thus allows for heterogeneity in shopping behavior. Our approach also allows a flexible model of shopping list that accounts for potential substitution within a good category, which can take place inside a store.

The main results of our empirical analysis indicate that the store-level own elasticities of expected basket cost are significantly biased toward zero when heterogeneity in price expectation is ignored. We also find that defining a shopping list at the product level somewhat mitigates the biases that stem from ignoring heterogeneity in price expectation but the own price elasticities are still substantially biased compared to the case of the basket-level shopping list.

This paper is, to the best of our knowledge, the first attempt to evaluate the effects of heterogeneity in price expectation on store choice. There is extensive marketing literature on price perception and reference prices in the context of brand choice. Most of this literature, starting from Monroe (1973) and Winer (1986), suggests that price information gathered from

various sources, combined with previously observed or paid prices, is integrated into consumers' expected or reference prices. More recently, Erdem, Imai, and Keane (2003) study the role of expectation on future prices in brand and quantity choices of forward-looking consumers. On the other hand, relatively little has been studied on the effects of price expectation in store choice decisions (see the related discussion in Mazumdar, Raj, and Sinha, 2005).

Our empirical work in this paper contributes to the literature on store choice problems in the supermarket industry. Smith (2004) examines merger effects with the estimated substitution patterns in a discrete and continuous choice setting and Beresteanu, Ellickson, and Misra (2010) estimate store choices using market level data and study the welfare effects of entry and competition in a dynamic framework. More recently, Katz (2007) estimates a store choice model using a moment inequality approach with household level data. The main difference between his work and ours lies on how to deal with the bundle-choice behavior in store choice. The moment inequality approach circumvents the need for dealing with the aforementioned complexity of choice sets by subtracting out all basket-related terms in a utility function except for those of interest. In the current study, we address the same problem by specifying the shopping basket composition in the first stage of the estimation.

More recently, Mojir, Sudhir, and Khwaja (2014) develop a model of spatiotemporal search that consumers can search across stores and across time to find the best possible prices. Their empirical analysis on store visits and purchases in the milk category shows that omitting the temporal dimension of search underestimates price elasticity (see also Gauri, Sudhir, and Talukdar (2008) for the relevance of consumers' search behaviors along the space and time dimensions). Although our setting does not explicitly model a temporal search, our use of duration model to predict consumers' purchasing time patterns can potentially reduce biases that may arise from omitting such temporal search behaviors.

# References

AILAWADI, K. L., AND S. A. NESLIN (1998): "The Effect of Promotion on Consumption: Buying More and Consuming It Faster," *Journal of Marketing Research*, 35(3), 390–398.

ALBA, J. W., S. M. BRONIARCZYK, T. A. SHIMP, AND J. E. URBANY (1994): "The Influence of Prior Beliefs, Frequency Cues, and Magnitude Cues on Consumers' Perceptions of Comparative Price Data," *Journal of Consumer Research*, 21(2), 219–235.

BAJARI, P., AND C. L. BENKARD (2003): "Discrete Choice Models as Structural Models of Demand: Some Economic Implications of Common Approaches," Working Paper.

BELL, D. R., D. CORSTEN, AND G. KNOX (2011): "From Point of Purchase to Path to Purchase: How Preshopping Factors Drive Unplanned Buying," *Journal of Marketing*, 75(1), 31–45.

BELL, D. R., T. H. HO, AND C. S. TANG (1998): "Determining Where to Shop: Fixed and Variable Costs of Shopping," *Journal of Marketing Research*, 35(3), 352–369.

BERESTEANU, A., P. B. ELLICKSON, AND S. MISRA (2010): "The Dynamics of Retail Oligopoly," Working Paper, Duke University.

BERRY, S., J. LEVINSOHN, AND A. PAKES (1995): "Automobile Prices in Market Equilibrium," *Econometrica*, 63(4), 841–890.

——— (2004): "Differentiated Products Demand Systems from a Combination of Micro and Macro Data: The New Car Market," *Journal of Political Economy*, 112(1), 68–105.

BLATTBERG, R. C., R. BRIESCH, AND E. J. FOX (1995): "How Promotions Work," *Marketing Science*, 14(3), G122–G132.

BLOCK, L. G., AND V. G. MORWITZ (1999): "Shopping Lists as an External Memory Aid for Grocery Shopping: Influences on List Writing and List Fulfillment," *Journal of Consumer Psychology*, 8(4), 343–375.

BODAPATI, A. V., AND V. SRINIVASAN (2006): "The Impact of Feature Advertising on Customer Store Choice," Mimeo, Stanford University.

BRIESCH, R., W. DILLON, AND E. FOX (2013): "Category Positioning and Store Choice: The Role of Destination Categories," *Marketing Science*, 32(3), 488–509.

BRIESCH, R. A., P. K. CHINTAGUNTA, AND E. J. FOX (2009): "How Does Assortment Affect Grocery Store Choice?," *Journal of Marketing Research*, 46(2), 176–189.

BRONNENBERG, B. J., M. W. KRUGER, AND C. F. MELA (2008): "Database Paper: The IRI Marketing Data Set," *Marketing Science*, 27(4), 745–748.

COX, D. R. (1972): "Regression Models and Life-Tables," *Journal of the Royal Statistical Society*, 34(2), 187–220.

DAVIS, P. (2006): "Spatial Competition in Retail Markets: Movie Theaters," *RAND Journal of Economics*, 37(4), 964–982.

ELLICKSON, P. B., AND S. MISRA (2008): "Supermarket Pricing Strategies," *Marketing Science*, 27(5), 811–828.

ERDEM, T., S. IMAI, AND M. P. KEANE (2003): "Brand and Quantity Choice Dynamics Under Price Uncertainty," *Quantitative Marketing and Economics*, 1(1), 5–64.

GAURI, K., K. SUDHIR, AND D. TALUKDAR (2008): "The Temporal and Spatial Dimensions of Price Search: Insights from Matching Household Survey and Purchase Data," *Journal of Marketing Research*, 45(2), 226–240.

GOOLSBEE, A., AND A. PETRIN (2004): "The Consumer Gains from Direct Broadcast Satellites and the Competition with Cable TV," *Econometrica*, 72(2), 351–381.

HANSEN, K., AND V. SINGH (2009): "Market Structure Across Retail Formats," *Marketing Science*, 28(4), 656–673.

HAUSMAN, J. A., AND D. A. WISE (1978): "A Conditional Probit Model for Qualitative Choice: Discrete Decisions Recognizing Interdependence and Heterogeneous Preferences," *Econometrica*, 46(2), 403–426.

HECKMAN, J. J., AND B. SINGER (1986): "Econometric Analysis of Longitudinal Data," in *Handbook of Econometrics*, ed. by Z. Griliches, and M. Intriligator, vol. 3. Amsterdam: North-Holland.

HENDEL, I., AND A. NEVO (2006): "Measuring the Implications of Sales and Consumer Inventory Behavior," *Econometrica*, 74(6), 1637–1673.

HO, T. H., C. S. TANG, AND D. R. BELL (1998): "Rational Shopping Behavior and the Option Value of Variable Pricing," *Management Science*, 44(12), S145–S160.

HOUDE, J.-F. (2012): "Spatial Differentiation and Vertical Mergers in Retail Markets for Gasoline," *American Economic Review*, 102(5), 2147–2182.

KALBFLEISCH, J., AND R. L. PRENTICE (2002): *The Statistical Analysis of Failure Time Data*. Hoboken, NJ: John Wiley.

KALWANI, M. U., AND C. K. YIM (1992): "Consumer Price and Promotion Expectations: An Experimental Study," *Journal of Marketing Research*, 29(1), 90–100.

KATZ, M. (2007): "Estimating Supermarket Choice Using Moment Inequalities," Unpublished Ph.D. dissertation, Harvard University.

KOLLAT, D. T., AND R. P. WILLETT (1967): "Customer Impulse Purchasing Behavior," *Journal of Marketing Research*, 4(1), 21–31.

KUMAR, V., AND R. P. LEONE (1988): "Measuring the Effect of Retail Store Promotions on Brand and Store Substitution," *Journal of Marketing Research*, 25(2), 178–185.

LAL, R., AND R. RAO (1997): "Supermarket Competition: The Case of Every Day Low Pricing," *Marketing Science*, 16(1), 60–80.

LANCASTER, K. (1971): *Consumer Demand: A New Approach.* New York: Columbia University Press.

LANCASTER, T. (1990): *The Econometric Analysis of Transitional Data.* Cambridge, UK: Cambridge University Press.

MAZUMDAR, T., S. P. RAJ, AND I. SINHA (2005): "Reference Price Research: Review and Propositions," *Journal of Marketing*, 69(4), 84–102.

MCFADDEN, D. (1981): "Structural Discrete Probability Models Derived From Theories of Choice," *Structural Analysis of Discrete Data and Econometric Applications.*

MOJIR, N., K. SUDHIR, AND A. KHWAJA (2014): "Spatiotemporal Search," COWLES Foundation Discussion Paper NO. 1942.

MONROE, K. B. (1973): "Buyer's Subjective Perceptions of Price," *Journal of Marketing Research*, 10(1), 70–80.

NEVO, A. (2001): "Measuring Market Power in the Ready-to-Eat Cereal Industry," *Econometrica*, 69(2), 307–342.

ORHUN, Y. (2013): "Spatial differentiation in the supermarket industry: The role of common information," *Quantitative Marketing and Economics*, 11(1), 3–37.

PAKES, A. (2010): "Alternative Models for Moment Inequalities," *Econometrica*, 78(6), 1783–1822.

PAUWELS, K., D. M. HANSSENS, AND S. SIDDARTH (2002): "The Long-Term effects of Price Promotions on Category Incidence, Brand Choice, and Purchase Quantity," *Journal of Marketing Research*, 39(4), 421–439.

PETRIN, A. (2002): "Quantifying the Benefits of New Products: The Case of Minivan," *Journal of Political Economy*, 110(4), 705–729.

SMITH, H. (2004): "Supermarket Choice and Supermarket Competition in Market Equilibrium," *Review of Economic Studies*, 71(1), 235–263.

THOMADSEN, R. (2005): "The Effect of Ownership Structure on Prices in Geographically Differentiated Industries," *RAND Journal of Economics*, 36(4), 908–929.

TRAIN, K. E. (2009): *Discrete Choice Methods with Simulation.* Cambridge, UK: Cambridge University Press.

WINER, R. S. (1986): "A Reference Price Model of Brand Choice for Frequently Purchased Products," *Journal of Consumer Research*, 13(2), 250–256.

# Appendix

## A  Alternative utility specification

To check the robustness of the (mean) utility specification we use for our estimation, we consider an alternative utility specification that includes a measure of assortment size and estimate the store choice, based on the utility

$$U_{ijbt} = \beta_i Dist_{ij} - \alpha_i Exp_{ijbt} + \delta_i Asrt_{ijbt} + \gamma_{ij} + \varepsilon_{ijt},$$

where $Asrt_{ijbt}$ denotes the variety of products carried by the store for the basket goods. The product variety in the model depends on the shopping list because households only consider the goods they plan to purchase. Thus, $Asrt_{ijbt}$ is defined by the average of the log of the number of products (or UPCs) carried by the store for the basket goods. Households may rather value the *overall* product variety for all goods instead of the basket goods. But we found that the overall assortment size independent of a shopping list is not statistically significant. This is because the overall assortment size barely changes over time at each store and it is rather absorbed by the store fixed effects in the model.

Although including a variety of available choices in a utility function is not a standard approach in consumer choice problems, allowing for good-specific product variety may be regarded as a remedy for abstracting away from the complicated brand choices in store choice estimation. However, as shown in Table 10, the own elasticities of expected basket cost only slightly change in magnitude without qualitative difference when product variety is included in the utility function.

## B  Explanatory variables in the duration models

We first construct a data-driven consumption variable for each good and then define inventory based on consumption levels. Consumption and inventory are defined at the goods level.[30] The underlying assumption in constructing consumption variables is that each household consumes

---

[30]Consumers may consume different brands for the same good and manage inventory at the brand level rather than at the goods level. Our setting does not allow brand-level inventories, and assumes good-specific inventory.

at a constant rate (or amount) for a set length of time. The amount of consumption may depend on various aspects of consumer heterogeneity, such as preference or special life events. Since none of these heterogeneous shocks is observable, we use a rather parsimonious approach to estimating consumption rates based on purchase frequency and quantity. Specifically, the constant consumption rate is computed as an average of the quantity of the products for the good purchased over a set period of time. Consumption is inferred for each household at each time, so it captures both time effects and heterogeneity of individual consumers.[31] Purchase quantity and frequency would be both associated with complex consumer behaviors, such as stock-piling and responses to promotional activities. Potential errors that can arise from neglecting to consider such consumer behavior in inferring the consumption level may be mitigated by averaging and smoothing the consumption rates over a length of time specified for each good.

The length of time for averaging the consumption level is set differently for each good, considering good-specific features such as average purchase frequency, shelf life, and whether it is a necessity good. For example, consumption rates are averaged over a short period of time for goods that are perishable and purchased with a high frequency such as milk and yogurt. On the other hand, for the goods that are storable (possibly with low-frequency of purchases) and that people typically consume in a persistent and constant pattern (e.g., necessity goods) we smooth consumption over a relatively long period (e.g., tissue, laundry detergent, and toothpaste).

Given the consumption level of each good, the inventory level at the beginning of each period is the residual amount in storage defined as follows:

$$I_{it} = \max\{0, \ I_{i,t-1} + Q_{i,t-1} - C_{i,t-1}\}$$

where $I_{it}$ is the inventory level consumer $i$ faces at time $t$, $Q_{it}$ is the quantity of the purchase at $t$, and $C_{it}$ is the estimated consumption level. The purchases during the first six months in the raw data are used to generate the distribution of initial inventories.

---

[31] Ailawadi and Neslin (1998) also construct time-varying consumption rates as a (continuous and nonlinear) function of inventory in their empirical study. Their approach heavily depends on functional specification.

Table 1. Demographics of Sample Customers

|  | IRI Data (%) | US Census (%) |
|---|---|---|
| Household income |  |  |
| 0–10,000 | 5.3 | 12.0 |
| 10,000–20,000 | 13.0 | 13.4 |
| 20,000–35,000 | 23.9 | 15.7 |
| 35,000–45,000 | 13.7 | 11.0 |
| 45,000–75,000 | 27.3 | 19.4 |
| 75,000– | 16.9 | 28.5 |
| Age of household head |  |  |
| <=25 | 1.2 | 7.1 |
| 25–44 | 20.7 | 29.3 |
| 45–54 | 24.5 | 20.9 |
| 55–64 | 21.7 | 16.1 |
| 65– | 32.0 | 26.5 |
| Single male | 3.6 | 23.7 |
| Single female | 28.2 | 27.6 |
| Home owners | 80.0 | 60.7 |
| Female employment | 38.6 | 51.5 |

*Notes:* The values in the table are the percentage of the sample customers or the population of the area. The population statistics are from the U.S. Census (The American Community Survey, 2005-2007)

Table 2. Summary Statistics of Household-Level Data

|  | Mean | Median | Std | Min | Max |
|---|---|---|---|---|---|
| Number of trips per month | 4.9 | 5 | 2.4 | 1 | 23 |
| Number of visited stores | 4.2 | 4 | 1.6 | 1 | 7 |
| Store HHI (visits) | 0.48 | 0.42 | 0.22 | 0.16 | 1 |
| Store HHI (spending) | 0.56 | 0.50 | 0.23 | 0.17 | 1 |
| Distance to store (mile) | 3.2 | 2.9 | 2.1 | 0.01 | 20.8 |
| Weekly spending (30 goods, $) | 20.83 | 15.54 | 18.71 | 0.16 | 363.16 |
| Weekly spending (all goods, $) | 102.81 | 84.20 | 81.23 | 0.13 | 1290.72 |
| Number of purchased goods per trip | 3.1 | 2 | 2.3 | 1 | 22 |

*Notes:* Store HHI of household $i$ is the sum of the square of the share of visits to (or dollar spending in) each store. That is, $StoreHHI_i = \sum_s \left(y_{its}/\sum_{s'} y_{its'}\right)^2$, where $y_{its}$ is the number of visits (or expenditure) of household $i$ for store $s$ at week $t$.



Figure 1. Store HHI by Visits and Spending

Figure 2. Distance and Store Choice

Table 3. Summary Statistics of Store-Level Data

|  | Chain | Price promotion (> 5 percent, %) | Assortment size | Average driving distance (mile) |
|---|---|---|---|---|
| Store 1 | A | 19.7 | 2,346 | 3.3 |
| Store 2 | A | 16.9 | 3,416 | 2.9 |
| Store 3 | B | 30.7 | 6,177 | 4.2 |
| Store 4 | B | 27.0 | 6,309 | 4.7 |
| Store 5 | C | 23.0 | 4,890 | 5.6 |
| Store 6 | C | 24.3 | 5,902 | 4.0 |
| Store 7 | D | 22.0 | 5,598 | 2.9 |

Figure 3. Price Expectation and Consumer Heterogeneity

Table 4. Estimation of Good Purchase Incidence

| | Blades | Soft drink | Cereal | Laundry detergent | Milk | Paper towel |
|---|---|---|---|---|---|---|
| Inventory (log) | -0.184*** | -0.642*** | -1.380*** | -0.127*** | -0.715*** | -0.229*** |
| | (0.054) | (0.005) | (0.009) | (0.015) | (0.006) | (0.013) |
| Consumption rate (log) | 0.644*** | 0.560*** | 0.673*** | 0.637*** | 0.655*** | 0.611*** |
| | (0.060) | (0.003) | (0.004) | (0.016) | (0.003) | (0.012) |
| Trip frequency | 0.029*** | 0.018*** | 0.004*** | 0.024*** | 0.011*** | 0.026*** |
| | (0.009) | (0.001) | (0.001) | (0.002) | (0.000) | (0.001) |
| Trip dollars (log) | 0.442*** | 0.112*** | 0.163*** | 0.294*** | 0.030*** | 0.298*** |
| | (0.088) | (0.005) | (0.001) | (0.017) | (0.004) | (0.014) |
| Weekend shopping | -0.002** | 0.000*** | 0.000 | 0.000 | 0.000* | 0.000 |
| | (0.001) | (0.000) | (0.000) | (0.000) | (0.000) | (0.001) |
| July 4th | -0.202 | 0.072*** | -0.094*** | -0.146*** | -0.020** | -0.072** |
| | (0.158) | (0.011) | (0.016) | (0.034) | (0.008) | (0.025) |
| Thanksgiving | -0.181 | 0.088*** | -0.122*** | -0.109** | 0.072*** | -0.101*** |
| | (0.179) | (0.012) | (0.020) | (0.038) | (0.008) | (0.029) |
| Christmas | 0.247 | 0.133*** | -0.107*** | -0.232*** | 0.070*** | -0.010 |
| | (0.174) | (0.011) | (0.019) | (0.039) | (0.008) | (0.028) |
| Purchase intervals | 0.009* | 0.002 | 0.006*** | 0.002 | -0.002 | 0.003* |
| | (0.004) | (0.002) | (0.001) | (0.002) | (0.002) | (0.001) |
| Volume size of purchases | 0.367*** | -0.104*** | -0.116*** | 0.026 | -0.201*** | 0.126*** |
| | (0.095) | (0.005) | (0.009) | (0.020) | (0.005) | (0.014) |
| Unit price of purchases | 0.091*** | 0.001 | 0.039*** | 0.026*** | 0.018*** | 0.018*** |
| | (0.013) | (0.003) | (0.004) | (0.003) | (0.002) | (0.003) |
| Pseudo R-squared | 0.19 | 0.04 | 0.10 | 0.11 | 0.03 | 0.07 |

*Notes:* Standard errors are in parenthesis (*$p < 0.05$, **$p < 0.01$, ***$p < 0.001$). Household demographics and purchases of other goods are also included in the estimation, but not reported here. Demographic variables included are log income, family size, and indicators for marriage, pet ownership and house ownership. Inventory level and volume size of purchases are normalized to the average quantity.
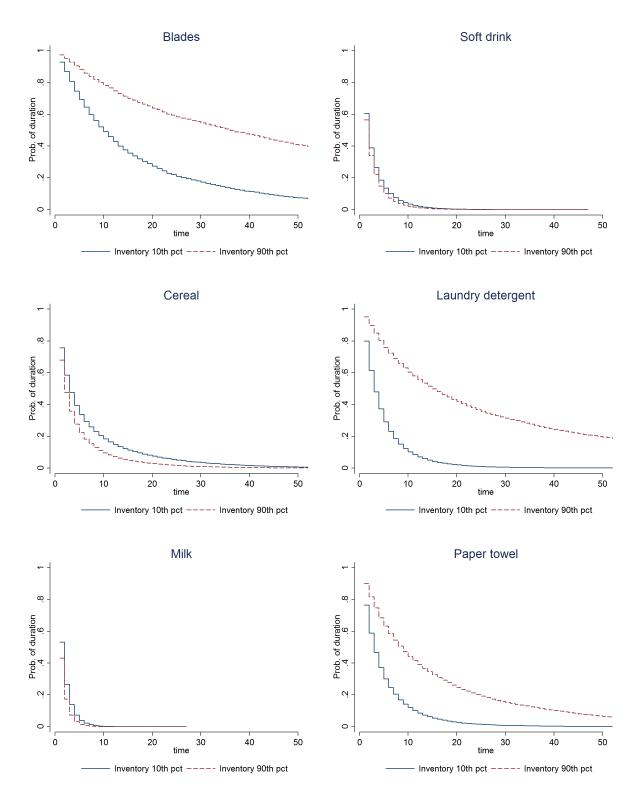
Figure 4. Inventory and Survivor Function

Table 5. Estimation of Quantity Choice

|  | Blades | Soft drink | Cereal | Laundry detergent | Milk | Paper towels |
|---|---|---|---|---|---|---|
| Price | -0.90*** | -16.45*** | -2.38*** | -2.96*** | -16.04*** | -32.72*** |
|  | (0.05) | (0.15) | (0.02) | (0.09) | (0.07) | (0.26) |
| Consumption | 0.42*** | 0.06*** | 0.04*** | 0.29*** | 0.10*** | 0.46*** |
|  | (0.07) | (0.00) | (0.00) | (0.01) | (0.00) | (0.01) |
| Trip freq. (good) | -0.09*** | -0.43*** | -0.02*** | -1.78*** | -0.18*** | -0.63*** |
|  | (0.02) | (0.02) | (0.00) | (0.08) | (0.01) | (0.02) |
| Trip freq. (overall) | -0.19 | -3.85*** | -0.08 | 2.38* | -3.53*** | 0.26 |
|  | (0.22) | (0.69) | (0.08) | (1.22) | (0.40) | (0.42) |
| R-squared | 0.41 | 0.16 | 0.22 | 0.23 | 0.36 | 0.40 |

*Notes:* Standard errors are in parenthesis (*$p < 0.05$, **$p < 0.01$, ***$p < 0.001$). Price is instrumented by the average price of the product in the same chain stores located in nine nearby cities. Fixed effects for time (year and month), brand, and customer are included.

Table 6. Errors in Expected Bundle Costs

|  | Expectation errors | | Actual |
|---|---|---|---|
| Frequency of trips | Full model | Non-heterogeneity | basket costs |
| 0< Frequency < 5 | 5.37 | 8.70 | 15.39 |
| 5≤ Frequency < 10 | 5.04 | 8.06 | 14.28 |
| 10≤ Frequency < 25 | 4.27 | 7.40 | 13.41 |
| 25≤ Frequency < 40 | 4.79 | 9.54 | 14.25 |
| 40≤ Frequency | 4.40 | 9.45 | 18.29 |

*Notes:* Expectation error is defined as the root mean square of the difference between expected cost and actual spending for each shopping basket. The full model reflects the sources of heterogeneity in price expectation, whereas non-heterogeneity model does not. Frequency is defined as the number of store trips per year.

Table 7. Store Choice Estimation: Full Model

| | | | Interactions with demographics | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Variable | Mean | Standard deviation | Income (15k-25k) | Income (25k-45k) | Income (45k-65k) | Income (65k+) | Age $\leq$34 | Age 35-54 | Age 55-64 | Child | Family size |
| Expected basket cost | -0.158 | 0.103 | 0.071 | 0.058 | 0.018 | 0.010 | 0.104 | 0.040 | 0.049 | – | – |
| | (0.017) | (0.007) | (0.020) | (0.019) | (0.020) | (0.019) | (0.023) | (0.012) | (0.013) | | |
| Distance | -0.441 | 0.443 | -0.094 | -0.061 | -0.162 | -0.133 | -0.027 | 0.057 | 0.016 | 0.003 | – |
| | (0.015) | (0.004) | (0.016) | (0.015) | (0.016) | (0.016) | (0.017) | (0.010) | (0.010) | (0.011) | |
| Store 2 (Chain A) | 0.257 | 1.166 | -0.384 | -0.487 | -0.363 | -0.317 | – | – | – | 0.357 | 0.031 |
| | (0.068) | (0.017) | (0.074) | (0.073) | (0.077) | (0.080) | | | | (0.063) | (0.021) |
| Store 3 (Chain B) | 0.802 | 1.155 | 0.483 | 0.339 | 0.691 | 0.700 | – | – | – | -0.274 | -0.056 |
| | (0.072) | (0.015) | (0.075) | (0.072) | (0.078) | (0.082) | | | | (0.059) | (0.021) |
| Store 4 (Chain B) | 0.549 | 1.107 | 0.029 | 0.233 | 0.664 | 1.074 | – | – | – | -0.260 | -0.077 |
| | (0.077) | (0.015) | (0.083) | (0.081) | (0.084) | (0.087) | | | | (0.067) | (0.024) |
| Store 5 (Chain C) | 0.454 | 1.336 | 0.603 | 0.498 | 1.043 | 1.001 | – | – | – | -0.245 | -0.109 |
| | (0.076) | (0.017) | (0.087) | (0.087) | (0.090) | (0.094) | | | | (0.069) | (0.024) |
| Store 6 (Chain C) | 1.502 | 1.169 | 0.041 | -0.061 | 0.453 | 0.404 | – | – | – | 0.108 | 0.022 |
| | (0.062) | (0.013) | (0.070) | (0.067) | (0.072) | (0.075) | | | | (0.056) | (0.018) |
| Store 7 (Chain D) | 0.705 | 1.016 | -0.221 | -0.108 | 0.212 | 0.255 | – | – | – | 0.185 | 0.051 |
| | (0.055) | (0.012) | (0.064) | (0.063) | (0.067) | (0.069) | | | | (0.057) | (0.019) |

*Notes*: The base group for age and income is older than 65 and earns less than 15 thousand dollars, respectively. Standard errors are computed by bootstrapping and given in parenthesis.

Table 8. Estimation of Store Choice: Model Comparison

| | Choice-based price expectation | | | Price expectation without heterogeneity | | |
|---|---|---|---|---|---|---|
| | Full model | No purchase quantity | No goods purchase | No shopping history | No choice-set | Product level basket |
| Expected basket cost | -0.158 | -0.130 | -0.165 | 0.089 | 0.153 | -0.063 |
| | (0.017) | (0.008) | (0.016) | (0.018) | (0.010) | (0.021) |
| Income | | | | | | |
| 15k-25k | 0.071 | -0.011 | 0.078 | -0.009 | -0.032 | 0.009 |
| | (0.020) | (0.008) | (0.019) | (0.021) | (0.013) | (0.025) |
| 25k-45k | 0.058 | 0.006 | 0.051 | -0.016 | -0.008 | 0.062 |
| | (0.019) | (0.008) | (0.017) | (0.018) | (0.011) | (0.022) |
| 45k-65k | 0.018 | -0.008 | 0.014 | -0.020 | -0.041 | 0.065 |
| | (0.020) | (0.009) | (0.018) | (0.020) | (0.012) | (0.021) |
| 65k+ | 0.010 | -0.013 | 0.036 | 0.009 | -0.057 | 0.097 |
| | (0.019) | (0.009) | (0.017) | (0.019) | (0.012) | (0.022) |
| Age | | | | | | |
| ≤34 | 0.104 | 0.053 | 0.114 | -0.003 | -0.007 | -0.047 |
| | (0.023) | (0.010) | (0.021) | (0.026) | (0.017) | (0.019) |
| 35-54 | 0.040 | 0.066 | 0.072 | 0.022 | 0.016 | -0.035 |
| | (0.012) | (0.005) | (0.011) | (0.009) | (0.007) | (0.009) |
| 55-64 | 0.049 | 0.036 | 0.047 | 0.012 | 0.013 | -0.089 |
| | (0.013) | (0.006) | (0.012) | (0.011) | (0.008) | (0.012) |
| Distance | -0.441 | -0.464 | -0.458 | -0.553 | -0.525 | -0.597 |
| | (0.015) | (0.012) | (0.013) | (0.013) | (0.017) | (0.014) |
| Income | | | | | | |
| 15k-25k | -0.094 | -0.055 | -0.098 | 0.145 | -0.023 | -0.015 |
| | (0.016) | (0.015) | (0.015) | (0.012) | (0.018) | (0.017) |
| 25k-45k | -0.061 | -0.004 | -0.022 | 0.138 | -0.022 | -0.017 |
| | (0.015) | (0.014) | (0.015) | (0.013) | (0.018) | (0.015) |
| 45k-65k | -0.162 | -0.101 | -0.047 | 0.108 | 0.009 | 0.113 |
| | (0.016) | (0.015) | (0.015) | (0.012) | (0.019) | (0.016) |
| 65k+ | -0.133 | -0.134 | -0.087 | 0.067 | 0.028 | 0.016 |
| | (0.016) | (0.015) | (0.021) | (0.013) | (0.019) | (0.017) |
| Age | | | | | | |
| ≤34 | -0.027 | -0.055 | 0.179 | 0.198 | 0.225 | 0.317 |
| | (0.017) | (0.016) | (0.019) | (0.021) | (0.020) | (0.028) |
| 35-54 | 0.057 | 0.030 | 0.119 | -0.022 | 0.170 | 0.213 |
| | (0.010) | (0.009) | (0.012) | (0.014) | (0.011) | (0.011) |
| 55-64 | 0.016 | -0.055 | -0.034 | -0.074 | 0.104 | 0.033 |
| | (0.010) | (0.011) | (0.010) | (0.014) | (0.012) | (0.011) |
| Child | 0.003 | 0.036 | 0.022 | -0.025 | 0.021 | -0.045 |
| | (0.011) | (0.011) | (0.012) | (0.011) | (0.011) | (0.012) |

*Notes:* The table only reports the estimates of the mean parameters and the interactions with demographics in the random coefficient model. Store fixed effects and interactions with demographics are included in all the estimations. Standard errors computed by bootstrapping are given in parenthesis.

Table 9. Own Price Elasticities

| | Choice-based price expectation | | | Price expectation without heterogeneity | | |
|---|---|---|---|---|---|---|
| | Full model | No purchase quantity | No goods purchase | No shopping history | No choice-set | Product level basket |
| S1-A | -5.43 | -7.39 | -6.17 | 4.59 | 4.67 | -2.33 |
| S2-A | -4.70 | -9.06 | -2.25 | 1.07 | 1.86 | -0.17 |
| S3-B | -4.79 | -5.50 | -6.67 | 2.11 | 7.59 | -2.95 |
| S4-B | -4.90 | -7.03 | -11.21 | 3.98 | 6.98 | -1.20 |
| S5-C | -1.91 | -2.59 | -1.34 | 1.30 | 1.84 | -3.18 |
| S6-C | -0.41 | -0.89 | -1.68 | 0.41 | 1.66 | -0.36 |
| S7-D | -0.80 | -1.59 | -1.90 | 0.43 | 1.68 | -0.34 |
| Mean | -3.28 | -4.86 | -4.46 | 1.98 | 3.76 | -1.50 |

*Notes:* Price elasticities are calculated as the percentage changes in the predicted market shares of store choices when expected basket costs of the sample customers increase by one percent. The prediction of store choices are based on the estimates from the random coefficient model.

Table 10. Robustness Check: Own Price Elasticities

| | Choice-based price expectation (Full model) | |
|---|---|---|
| | With assortment size | Without assortment size |
| S1-A | -7.07 | -5.43 |
| S2-A | -6.61 | -4.70 |
| S3-B | -5.10 | -4.79 |
| S4-B | -3.37 | -4.90 |
| S5-C | -2.87 | -1.91 |
| S6-C | -0.77 | -0.41 |
| S7-D | -1.29 | -0.80 |
| Mean | -4.01 | -3.28 |