# Do Consumers Pay for Being Healthy Conscious?—An Analysis of Price Discrimination On Healthier Food Product

Congnan Zhan

April 25, 2010

North Carolina State University

Department of Economics

Campus Box8110, Raleigh

NC27695-8110

czhan@ncsu.edu

**Abstract**

'Healthier food product' has experienced a rapid growth rate in recent years in U.S. because of the increasing consumer demand for healthier and environmental friendlier lifestyle. This analysis is looking for price discrimination evidences by comparing price cost margins of regular food product and healthier food products. Price cost margins are computed by solving firms' profit maximization problem and relevant parameters are estimated from consumers' choice decisions. Specifically, price elasticities and price coefficients are estimated using nonlinear GMM estimation in order to construct price cost margin. The empirical analysis employs product level ketchup data across 50 MSA in U.S. from 2001 to 2006.

Key Words: Price Discrimination, Discrete Choice Models, Random Coefficients, Ketchup Industry.

# 1  Introduction

'Healthier food product' has experienced a rapid growth rate in recent years in U.S. because of the increasing consumer demand for healthier and environmental friendlier lifestyle. People's enthusiasm in health, food safety, environments and even animal welfare cause the increasing demand for food products such as organic food, reduced/low sugar food, reduced/low fat food and so on. Organic food guarantees no usage of antibiotics and hormones in livestock production and the use of organically grown feed and pasture. Price differences are observed between healthier food products and conventional food products. The observed price differences could be possibly due to different consumer willingness to pay, different production cost or producer price discrimination. This paper tries to answer the question that whether 'healthier food' is over-priced by producers and whether more healthy-oriented consumers are price-discriminated. Do they pay for being healthy conscious?

There are papers studying the production side of organic food, policy responses, producers' competition and consumers' willingness to pay for healthier food. But, no existing published study examines the existence of price discrimination of healthy-oriented food product. Villas-Boas and Zhao (2005) is a study which systematically model the ketchup industry. It develops both demand and supply sides of the market. It mainly focuses on manufacture competition and retailer-manufacture interactions. Several WTP studies try to answer the question if consumers are willing to pay more money to organic food and how much more. Among them, Batte Hooker Haab and Beaverson (2007) concludes that consumers are willing to pay premium price for organic even though products contain less than 100% organic ingredients.

Regarding methodology, the paper aims to recover and compare the price cost margin of individual organic and reduced sugar food products and conventional food products by estimating the demand system. Evidence of over pricing for organic/reduced sugar food will be obtained if the average markup of organic/reduced sugar food is significantly larger than that of the conventional food products. In more detail, the paper will apply the BLP (1995) and Nevo(2001) method to recover the price cost margin by estimating the demand system and solving the firm's profit maximization problem .

The analysis is on U.S. nation-wide product level using the scanner dataset from IRI and ketchup is the targeted food product. There are 50 markets in the dataset. Approximately,

the population ranges from 19,000,000 to 45,000 with average 3,450,000. The 50 markets are divided into 4 regions: North East, West, MidWest and South. Organic and reduced sugar ketchup products are chosen to represent 'healthier food product'. In this analysis, we will focus on 92 products (at the UPC level), which cover about 95% of the total market in six consecutive years. Among 92 products, 55 of them are produced by Heinz and Conagra food (Hunt is its famous brand), accounting for about 80% of the total market.

The paper is organized as follows. Section 2 provides a description about the ketchup industry. Section 3 describes the empirical framework of the model while section 4 is about estimation method, instruments used, UPC dummies and data. Section 5 provides results of Logit model. Section 6 concludes.

## 2  Ketchup Industry

The analysis assume two hypothetical firms are competing in the industry, firm one producing only regular products and firm two producing only healthier products. Price cost margins are able to be backed out for products in each firm. Ketchup industry is chosen because of its concentration of market structure. Ketchup is the most widely used condiment in the US in 97% of all kitchens. Currently 56% of ketchup is consumed with three main foods: hamburgers, hot dogs and french fries, which remain the most eaten foods.[1] Heinz is the largest ketchup producer. Hunt is the second largest brand followed by Del Monte, Generic and some private labels. The volume market shares in 6 years are summarized in Table1.

*Please Table1 here:*

According to the table above, the total market share of Heinz, Hunt's and Del Monte ranges from 78% to 80% from 2001 to 2006. Heinz itself accounts for about 60% of the total market share. The analysis therefore includes all products of Heinz, Hunt's and some other products that are organic, reduced sugar or no salt added.

---

[1]Survey of national eating trends by NPD Group

# 3  Empirical Framework

## 3.1  Derive price cost margin from frim's problem

The objective of the analysis is to estimate and compare the price cost margin between regular product and healthier product. The empirical framework follows BLP(1995) and Nevo(2001). Price cost margin is derived from firms profit maximization problem. Market share derivatives are constructed from parameters which are estimated from consumers' choice. Thus, firstly, demand of all products is estimated. Secondly, own and cross price elasticities are computed in order to obtain the price cost margin.

The model assumes two firms in the ketchup industry, one firm produces regular products and the other produces the special category. J products in total are produced. Each product owns a unique UPC. Firms profit maximization problem is as follows:

$$Max\Pi_f = \sum_{j \in \mathscr{F}_f} (p_j - mc_j)Ms_j(p) - C_f \tag{1}$$

$$M_k = Max\left(\sum_{week} \sum_{store} \sum_{UPC} q_{j,store,upc,k,week}\right) \tag{2}$$

where $\mathscr{F}_f$ is the subset of J and represents all UPCs belonged to firm $f$. $p_j$ and $mc_j$ are price in ounce and marginal cost of product j. M is the potential market size which is defined as the largest aggregate value of units purchased within 72 months in each market. $s_j(p)$ is the market share of product $j$. And $C_f$ is the fixed cost of production for firm $f$.

By taking the FOC, the following equation is obtained.

$$s_j(p) + \sum_{r \in \mathscr{F}_f} \left((p_r - mc_r)\frac{\partial s_r(p)}{\partial p_j}\right) = 0 \tag{3}$$

Define $pcm = (p_j - mc_j)$ for $\forall j$

$$s_j(p) = -\sum_{r \in \mathscr{F}_f} \left(pcm \frac{\partial s_r(p)}{\partial p_j}\right)$$

Define $S_{jr} = -\frac{\partial s_r(p)}{\partial p_j}$, then $s_j(p) = \Sigma(pcm * S_{jr})$

When $S_{jr}$ and $s_j(p)$ are known, $pcm$ are obtained by solving the linear equations.

Define $\Omega_{jr} = \Omega_{jr}^* * S_{jr}$, where $\Omega_{jr}^* = 1, r \in \mathscr{F}_f, j \in \mathscr{F}_f 0, otherwise$

Thus, $s(p) = \Omega * pcm$

$$pcm = \Omega^{-1} * s(p) \tag{4}$$

From supply side, price cost margin is derived in term of market shares and derivatives of market shares. These need to be estimated from demand side.

## 3.2 Consumer's Problem

Price elasticities of market share are derived from consumer's problem. The specification of consumer i's utility function $U(x_{jt}, p_{jt}, \xi_j, D_i; \theta)$ is a function of observed product characteristics $x$, unobserved product characteristics $\xi$, product price by ounce, demographic characteristics and unknown parameter $\theta$. Following [**?**], the specification is

$$u_{ijt} = x_{jt}\beta_i^* - p_{jt}\alpha_i^* + \xi_j + \triangle\xi_{jt} + \varepsilon_{ijt} \tag{5}$$

$$, i = 1, ...I, j = 1, ...J, t = 1, ..., T$$

where $x_{jt}$ is a K-dimensional vector containing observed product characteristics. $x_{jt}$ includes a healthy dummy and size of each product. $\xi_j$ is the mean value of unobserved product characteristics while $\triangle\xi_{jt}$ is the time-product deviation from the mean value. $\varepsilon_{ijt}$ is a mean zero stochastic term.

$[\alpha_i^*; \beta_i^*]$ are parameters describe choice of consumer i. It is a $K + 1$ dimension column vector. Consumers' preferences vary as a function of observed individual characteristics and unobserved characteristics. Let $[\alpha; \beta]$ be the mean value of $[\alpha_i^*; \beta_i^*]$. Therefore, the representation of $[\alpha_i^*; \beta_i^*]$ is

$$[\alpha_i^*; \beta_i^*] = [\alpha; \beta] + \Pi D_i + \Sigma v_i, v_i N(0, I_{K+1}) \tag{6}$$

where $D_i$ is the demographic variables including individual income, individual income square, education and household size. $v_i$ is unobserved consumer tastes which are random draws

from multi-variate normal distribution. $\Pi$ and $\Sigma$ are parameters need to be estimated. $\Pi$ includes all interaction term between observed product characteristics and observed demographic characteristics. $\Sigma$ is a scaling matrix.

By merging equation (5) and (6),

$$u_{ijt} = \delta_{jt}(p_{jt}, x_{jt}, \xi_{jt}, \triangle\xi_{jt}; \theta_1) + \mu_{ijt}(p_{jt}, x_{jt}, D_{it}, v_{it}; \theta_2) + \varepsilon_{ijt} \tag{7}$$

where $\delta_{jt}(p_{jt}, x_{jt}, \xi_{jt}, \triangle\xi_{jt}; \theta_1) = x_{jt}\beta - p_{jt}\alpha + \xi_j + \triangle\xi_{jt}$

and $\mu_{ijt}(p_{jt}, x_{jt}, D_{it}, v_{it}; \theta_2) = [-p_{jt}, x]$

$\delta_{jt}$ is the mean utility of product $j$ at time $t$ while $\mu_{ijt}$ is the deviation of each individual. $(\mu_{ijt} + \varepsilon_{ijt})$ is mean zero heteroskedasticity deviation from the mean.

In this discrete choice model, an outside good need to be define to include the case consumer $i$ does not choose any product. In another word, it is likely that the sampled individual is not a consumer of product $j$. The specification of outside good is as follows.

$$u_{i0t} = \xi_0 + \pi_0 D_i + \sigma_0 v_i + \varepsilon_{i0t} \tag{8}$$

The consumers' problem is choosing one unit of product that gives highest utility assuming no tie occurs. Mathematically,

$$A_{jt}(x, p._t, \delta._t; \theta_2) = \{(D_i, v_i, \varepsilon_{it}) | u_{ijt} > u_{ilt} \forall l = 1, ..., J\} \tag{9}$$

where $\theta_2 = [\Pi, \Sigma, \pi_0, \sigma_0]$

The predicted market share is defined as:

$$s_{jt}(p._t, x, \delta._t; \theta_2) = \int_{A_{jt}} dP(D_i, v_i, \varepsilon_i) = \int_{A_{jt}} dP(\varepsilon_i | D_i, v_i) dP(v_i | D_i) dP(D_i) = \int_{A_{jt}} dP(\varepsilon_i) dP(v_i) dP(D_i) \tag{10}$$

$$, v_i \sim N(0, I_{k+1})$$

$\varepsilon_i$ is the mean zero stochastic term. $P(.)$ is the distribution function which is approximate by

sampling CPS.

The model becomes multinomial Logit by simply assuming $\varepsilon_{ijt}$ is iid extreme value distribution and consumers heterogeneity only enter the model through separable additive random shock. Even though it provides a closed functional form for equation(10), it restrict the substitution pattern of own and cross price elasticities (Nevo(2001), Nevo(1995)). The market share and elasticities under Logit market are as follows:

$$s_j = \frac{exp(x_{jt}\beta^* - \alpha^* p_{jt} + \xi_j + \triangle \xi_{jt})}{1 + \sum_{k=1}^{J}(exp(x_{jt}\beta^* - \alpha^* p_{jt} + \xi_j + \triangle \xi_{jt}}$$  (11)

The price elasticities of market shares for Logit are:

$$e_{jkt} = \frac{\partial s_{jt}}{\partial p_{kt}} \frac{p_{kt}}{s_{jt}} = \begin{cases} \alpha^* p_{jt}(1 - s_{jt}), & j = k \\ -\alpha^* p_{jt} s_{jt} & , otherwise \end{cases}$$  (12)

It generates two major problems.[2] Firstly, the own price elasticity proportionally depends on the its own price so that low price products has smaller elasticities which implies a high markup. This is potentially problematic because it's possible when marginal costs of cheaper products are lower as percentage of price. Secondly, the substitution pattern is restricted by the form of market share. If the market share of two products from two exclusive categories are the same, when price of another product from either group increase, substitution from this product will toward both groups. The second argument is not a problem for this analysis because the products are segmented into healthier and others. For example, if price of organic ketchup increase, consumers in this group are likely substitute to regular ketchup and other organic product.

However, because of the first issue above, Logit is abandoned and the full model allows correlation between unobserved variables and reasonable substitution patterns. Instead of assuming iid extreme value $\varepsilon_{ijt}$, it assumes variance components structure. The market share of individual sampled consumer and elasticities of products are

$$s_j = \frac{exp(\delta_{jt} + \mu_{ijt})}{1 + \sum_{k=1}^{J} exp(\delta_{kt} + \mu_{ikt})}$$  (13)

---

[2]There is a detailed discussion in Nevo's "A Research Assistant's Guide to Random Coefficients Discrete Choice Models of Demand " 1998 technocal working paper.

and

$$\eta_{jkt} = \frac{\partial s_{jt}}{\partial p_{kt}} \frac{p_{kt}}{s_{jt}} = \begin{cases} (\frac{p_{kt}}{s_{jt}} \int \alpha_i s_{ijt}(1 - s_{ijt}) dP(D) dP(v)) & , j = k \\ (-\frac{p_{kt}}{s_{jt}} \int \alpha_i s_{ijt} s_{ikt}) dP(D) dP(v)) & , otherwise \end{cases} \qquad (14)$$

This solves the first issue above by a functional form of elasticity which doesn't only depend on its own price.

# 4  Data

The first dataset has 118897 observations including average ounce price, observed market share, healthy dummy, size, month dummy and UPC dummy etc. The data covers 50 U.S. regions ranging from 2001 January to 2006 December. The region definition of IRI marketing dataset is close to definition of MSA (Metropolitan Statistical Area) in CPS. In estimation, one market implies one region-month combination. For example, market one is Atlanta in 2001 January. Region Atlanta includes counties around Atlanta. Also, 50 U.S. regions are aggregated into 4 large regions, North East, Mid West, West and South. This aggregation is needed in order to construct instrumental variables.

The analysis is on UPC level. Ketchup products from same brand with different flavors have different UPC. This property of IRI data set provides a chance to identify product characteristics. Among 401 products, all Heinz and Hunt's that showed up at least once in 6 years are included. UPCs having annual market share bigger than 0.1 and in the market in all 6 years are included. Besides, all UPCs that are organic, reduced sugar or no sugar are also included no matter the value of market share. One problem about the UPC is some IRI recorded UPCs are different from UPC bar code found on product description label. In this case, a criterion is set up to merge products sharing same properties. The detail explanation is also provided in appendix A. In total, 90 UPCs (after merging) are included and the amount of UPC included in each market varieties. Under each UPC-market combination, price is constructed as dollars sales divided by quantity sold in ounces while ounces are recorded from product label. Average ounce price is adjusted inflation by CPI of food segment in each aggregate region.

Market share is defined as potential market share and constructed as:

$$S_{jm} = \frac{q_{jm}}{PotentialMarketSize} \forall m = 1,...,M$$

The data used to estimate demand parameters consists of price, observed market shares, product characteristics and demographic characteristics. These data come from two main sources, IRI Marketing Data Set and Current Population Survey from BLS. Price, market shares and product characteristics are constructed from IRI marketing dataset and demographic characteristics are sampled from BLS. Details on how these data are constructed are provided in appendix A. where M is the total number market and $q_{jm}$ is the quantity in ounce of product j sold in market $m$. $q_{jm}$ is aggregated from quantity sales in all stores in market $m$(A month-region combination). Potential market size is defined as the largest volume sold in one month summing up all UPCs in all stores in one region across 72 months. The outside good market share is difference 1.01 and sum of inside good shares in order to avoid zeros outside goods for those 72 markets which volume sum is the maximum in that region. Product characteristics are size and healthy dummy variable. Healthy dummy is one if observation is organic, reduced sugar, no sugar or no salt. The following table summarizes statistics of price and market share.

The second dataset includes demographic variables such as income, income square, education and age corresponding to each market-year combination. In each market, 20 individual are sampled from CPS. Table 2 summarizes demographic statistics.

*Place Table2 here:*

# 5   Econometrics

The estimation employs GMM estimator following Nevo(2001) and BLP (1994). The moment conditions are assumed to be zero such that,

$E[Z'\omega(\theta)] = 0$, where $\omega(\theta)$ is the error term and $Z$ includes instruments. In this discrete choice model, $\omega(\theta) = \xi_j + \triangle\xi_{jt}$, which is derived from the previous equation $\delta_{jt}(p_{jt}, x_{jt}, \xi_{jt}, \triangle\xi_{jt}; \theta_1) = x_{jt}\beta - p_{jt}\alpha + \xi_j + \triangle\xi_{jt}$.

Therefore,

$$\xi_j + \triangle\xi_{jt} = \delta_{jt}(p_{jt}, x_{jt}, \xi_{jt}, \triangle\xi_{jt}; \theta_1) - x_{jt}\beta - p_{jt}\alpha + \xi_j. \tag{15}$$

The GMM objective function is $\omega(\theta)'ZA^{-1}Z'\omega(\theta)$, where $A^{-1}$ is an optimal weight matrix constructed from instrumental variables. Specifically, $A^{-1} = (Z'Z)^{-1}$. By searching for $\theta$, the program is minimizing GMM objective function. In order to construct error therm $\omega(\theta)$, $\delta_{jt}$ needs to be available. Because $\delta_{jt}$ doesn't have a closed specification, it is obtained by numerically contraction mapping from following equation.

$$s_{.t}(x, p_{.t}, \delta_{.t}; \theta_2) = S_{.t}$$

, where $S_{.t}$ is the observed potential market share and $s_{.t}$ is the predicted market share. The starting value of $s_{.t}$ is the results of 2SLS. From the results of contraction mapping, parameters enter the model linearly ($\theta_1 = [\alpha, \beta]$) is obtained. Using the estimates of $\theta_1$, $\theta_2$ is estimated in GMM. The weight matrix A is computed in two steps. Firstly, optimal weight matrix is used to get estimates of parameters and secondly, using initial estimates A is computed again in order to reduce the variances. After obtaining $\theta_2$, own and cross elasticities are computed to obtain price cost margins.

The instrumental variables in the estimation include monthly regional average prices and cost proxy. The first set of instruments is constructed from prices. They are monthly average regional prices. Specifically, in each region, the prices of markets excluding observation market are averaged in each month. This average price is reported as an instrument for products of the observation market. When products are only sold in one market in one region, the average price of all other region in the month is used. Also, when products are not sold in some months, average prices across months are used as instruments. 20 regional prices are selected at the end. The second category of instruments are cost proxy including regional dummy, hourly wage in supermarket sector in each market each month and population density of each MSA. These instruments are also employed in both full model and Logit estimation.

In addition, UPC dummies are included in Logit and full model estimation to capture the true factors that determine utilities. The UPC dummies enter the model linearly and don't increase the estimation difficulty. In Logit estimation, when products observed characteristics are not included, UPC dummies are included. The results are discussed in following section.

# 6 Estimation Results

## 6.1 Logit Results

Even though the Logit model generates restrictive substitution pattern, the estimate are computed to test the importance of different combination of instrumental variables. Part of results of Logit regression is as in following table. The dependent variable is $ln(s_{jt}) - ln(s_{0t})$, which is the mean utility of product $jt$ in Logit model. OLS and IV regressions are performed to test the strength of different set of IV. Table 3 shows the results from Logit model.

*Place Table 3 here:*

The column i to iii are OLS regressions. Column i regresses on monthly dummies, product characteristics and price.Column ii additionally includes UPC dummies and iii includes demographic variables individual mean income, mean education and mean household size. The results in ii and iii are very close. All of the reported coefficients in OLS are significant. The column iv to viii are 2SLS using different IVs. Column iv uses UPC dummies and returns results close to OLS (i). Because in column i products characteristics are included and in iv UPC dummies also includes all of these variables. The results in column v and vi reports IV regression using monthly average regional prices as instruments. Also as column ii and iii, by adding demographic regressors the results do change much. Also average regional prices reduce the magnitudes of price coefficients. The demographic coefficients suggest that the mean utility of consuming ketchup increases with mean income while education coefficients are not significant. The last two columns are regressions using both prices and cost proxies as instruments. The $R^2$ of OLS and first stage $R^2$ of IV regressions are big except the first column. Also, the first stage F tests show big values indicating the strength of instruments.

# 7 Conclusion and Future Work

The paper apply discrete choice model to estimate demand coefficients of two groups of ketchup products. The results from GMM estimation are used to compute elasticities which are required to obtain price cost margins for both regular and special categories ketchup products. The results of Logit tested the efficiency of different set of instruments. It implies that average regional prices and UPC dummies work for model generating significant and realistic coefficients while

the cost proxies doesn't accurately predict the mean utility. Besides, all regressions show right relationship between mean utility and prices. These could serve as starting values in full model GMM estimation. The GMM algorithm is going to be applied to this analysis in the future based on the Logits results above. Furthermore, the price cost margins will be computed based on estimates of GMM.

# 8 Appendix

## 8.1 Appendix A: Introduce IRI data

The data required in the analysis are from IRI Marketing Dataset and BLS. IRI dataset provides weekly dollars sales, quantity sales and product UPC for each sampled store in each region. Regions are divided into south, west, midwest and north east by author. IRI data were collected using scanner devices in randomly selected super market and grocery stores in 50 regions across U.S. IRI region definition is close to MSA definition from BLS including metropolitan area and rural area towns surrounded. Price, market share, product characteristics are all provided by IRI data.

Market share is computed in the unit of ounces. It's the total ounces sales of one UPC in a given month in a given market derived by the total potential market size in the corresponding market. For all months in the same market, potential market size is the same which is defined as the largest volume sales in this market across 72 months. The outside good of each month-market combination is 1.01 subtracts the total market shares in given month-market combination. The price variable is the constructed by dividing dollars sales(IRI data) by total ounces sold for given UPC. The unit of price is cent per ounce sold of given UPC in given month-market combination without manufacture coupons. The variable "dollars" in IRI data is the retail price paid, on average including retail features, displays, and retailer coupons and excluding manufacturer coupons or any discount that might be applied by the retailer that is not applicable to the item. For example, if a retailer gave $5 off if you purchased more than $200, that discount is not applied.[3] Besides, price is adjusted inflation using seasonally adjusted monthly regional food CPI. 50 markets are categorized into 4 regions and each region is divided into group A

---

[3]This explanation is cited from " IRI academic dataset description"

and B according to the population size.The criteria to categorize population is 1,500,000. Both population and CPI data are from BLS.

Product characteristics are also provided by IRI academic dataset. It describes the size, sugar content and style of each UPC product. Organic is one of the styles. No sugar and low sugar are described in sugar content. No sugar, low sugar, organic and no salt are combined as healthy dummy. The size of bundled products is the individual product size times quantity bundled. For example, some UPC includes two bottles of ketchup with same size. The size for this UPC is individual bottle size times two.

Demographic variables are collected from BLS March Current Population Survey. Demographic variables include income, income square, education and household size. 20 individuals are sampled in each market each year. All of the three variables are household variables except education is individual variable. Income is defined as individual income which created by dividing household income by household size. The mean of sampled income is lower than BLS national average income because it doesn't exclude individual under 16 years old.

Additional two instrumental variables are employed in IV regression. Average hourly earnings in supermarket sector in each market and population density in each city are from NBER CPS Monthly Earning Extracts and BLS respectively. They serve as regional price indices in IV regression.

There are 401 UPC in IRI data from 2001 to 2006. The criterion to pick UPC is as follows. Firstly, all Heinz and Hunt's UPC that showed up at least once in 2001 to 2006. Secondly, other UPC annual market share>0.1 and show up in 6 year consistently. Lastly, UPC that are organic or reduced/no sugar and not from Heinz and Hunt's. There are 105 out of 401 UPC picked. One problem of UPC is that different IRI UPCs possibly represent one product. This is due to three reasons. One of the reason is IRI updated the generation code. Products with different generation code could be the same. Secondly, when bundled products are sold separately IRI reindex it as a UPC starting with system code 27(IRI UPC is converted from PLU and SKU). Therefore, products having different system code could represent same product. Lastly, missing data are found in style and sugar content variables. It's hard to distinguish the difference between products in same size, brand and company. A criterion is set up to remove duplicates in selected 105 UPCs. UPCs are counted as the same product as long as they are under same brand and

same producer and have same attributes in sugar content and regular. The relevant products are reindexed and number of UPC reduced to 90. Collapsing data as above will not significantly affect the results of analysis because the objective is to study the competition between healthy-oriented products (organic/reduced sugar) and regular products.

Furthermore, Table 4 provides a summary of all brand and number of UPC included in each category. And Table 5 summarizes demographic and product characteristic variables.

*Place Table 4 here:*

*Place Table 5 here:*


## 8.2   Appendix B: Derivation of Consumer's Problem

Consumers' Problem:

Thus, the utility specification includes all of observed and unobserved characteristics as follows.

$$u_{ijt} = x_j \beta_i^* - p_{jt} \alpha_i^* + \xi_j + \Delta \xi_{jt} + \varepsilon_{ijt} \tag{16}$$

where $j = 1, ...J; i = 1, ...I$

$I$ individual consumers are sampled in each market. $x_{jt}$ represents the product $j$'s characteristics in time $t$. Specifically, in this empirical analysis, $x = [constant \quad healthy \quad size]$. $p_{jt}$ is the ounce price of product $j$ in time $t$. $\xi_j$ is the mean value of unobserved characteristics of product $j$ while $\Delta \xi_{jt}$ is the deviation of time $t$ or market $t$ from mean value of product $j$. $\varepsilon_{ijt}$ has zero mean including all other shocks that are able to change utility.

The specification of outside good is as follows.

$$u_{i0t} = \xi_0 + \pi_0 D_i + \sigma_0 v_i + \varepsilon_{i0t} \tag{17}$$

The mean value of product characteristics is normalized to zero.

For different consumer $i$, the valuation of specific product characteristics varies according to various consumers' characteristics. For example, people who are more educated may tend to focus more on nutrition of a product while some other groups pay less attention on this. Also, consumers from large household size may choose big packaged products. Thus $\alpha$ and $\beta$

15

varies from i. The deviation includes observed demographic variables income, education and household size and unobserved variables $v_i$.

$$
\begin{bmatrix} \alpha_i^* \\ \beta_{i,constant}^* \\ \beta_{i,healthy}^* \\ \beta_{i,size}^* \end{bmatrix} = \begin{bmatrix} \alpha_i \\ \beta_{i,constant} \\ \beta_{i,healthy} \\ \beta_{i,size} \end{bmatrix} + \begin{bmatrix} \Pi_{11} & \cdots & \Pi_{1,4} \\ \vdots & \ddots & \vdots \\ \Pi_{4,1} & \cdots & \Pi_{4,4} \end{bmatrix} \begin{bmatrix} D_{i,income} \\ D_{i,income^2} \\ D_{i,education} \\ D_{i,hhsize} \end{bmatrix} + \begin{bmatrix} \Sigma_{11} & \cdots & \Sigma_{1,4} \\ \vdots & \ddots & \vdots \\ \Sigma_{4,1} & \cdots & \Sigma_{4,4} \end{bmatrix} \begin{bmatrix} v_{i,1} \\ v_{i,2} \\ v_{i,3} \\ v_{i,4} \end{bmatrix}
$$

(18)

$$
v_i \sim N(0, I_4)
$$

where, $\Pi$ is the parameter matrix explaining demographic varieties and $\Sigma$ is a scaling matrix. The first vector on the right hand side is the mean value of $\alpha$ and $\beta$.

Thus, by substituting (B3) into equation (B1), the individual utility function is:

$$
u_{ijt} = [-p_{jt}, x_{j,1}, x_{j,2}, x_{j,3}] \begin{bmatrix} \alpha_i \\ \beta_{i,constant} \\ \beta_{i,healthy} \\ \beta_{i,size} \end{bmatrix} +
$$

$$
[-p_{jt}, x_{j,1}, x_{j,2}, x_{j,3}] \left\{ \begin{bmatrix} \Pi_{11} & \cdots & \Pi_{1,4} \\ \vdots & \ddots & \vdots \\ \Pi_{4,1} & \cdots & \Pi_{4,4} \end{bmatrix} \begin{bmatrix} D_{i,income} \\ D_{i,income^2} \\ D_{i,education} \\ D_{i,hhsize} \end{bmatrix} + \begin{bmatrix} \Sigma_{11} & \cdots & \Sigma_{1,4} \\ \vdots & \ddots & \vdots \\ \Sigma_{4,1} & \cdots & \Sigma_{4,4} \end{bmatrix} \begin{bmatrix} v_{i,1} \\ v_{i,2} \\ v_{i,3} \\ v_{i,4} \end{bmatrix} \right\} + \xi_j + \Delta \xi_{jt} +
$$

$\varepsilon_{ijt}$,

rearranging,

$$
u_{ijt} = [-p_{jt}, x_{jt,1}, x_{jt,2}, x_{jt,3}] \begin{bmatrix} \alpha_i \\ \beta_{i,constant} \\ \beta_{i,healthy} \\ \beta_{i,size} \end{bmatrix} + \xi_j + \Delta \xi_{jt}
$$

$$
+ [-p_{jt}, x_{jt,1}, x_{jt,2}, x_{jt,3}] \left\{ \begin{bmatrix} \Pi_{11} & \cdots & \Pi_{1,4} \\ \vdots & \ddots & \vdots \\ \Pi_{4,1} & \cdots & \Pi_{4,4} \end{bmatrix} \begin{bmatrix} D_{i,income} \\ D_{i,income^2} \\ D_{i,education} \\ D_{i,hhsize} \end{bmatrix} + \begin{bmatrix} \Sigma_{11} & \cdots & \Sigma_{1,4} \\ \vdots & \ddots & \vdots \\ \Sigma_{4,1} & \cdots & \Sigma_{4,4} \end{bmatrix} \begin{bmatrix} v_{i,1} \\ v_{i,2} \\ v_{i,3} \\ v_{i,4} \end{bmatrix} \right\} + \varepsilon_{ijt},
$$

then,

$$
u_{ijt} = \delta_{jt}(p_j, x_j, \xi_j, \Delta \xi_{jt}; \alpha_i, \beta_i) + \mu_{ijt}(p_j, x_j, D_i, v_i; \Pi, \Sigma) + \varepsilon_{ijt},
$$

where,

$$\delta_{jt}(p_j,x_j,\xi_j,\Delta\xi_{jt};\alpha_i,\beta_i) = [-p_{jt},x_{j,1},x_{j,2},x_{j,3}]\begin{bmatrix} \alpha_i \\ \beta_{i,constant} \\ \beta_{i,healthy} \\ \beta_{i,size} \end{bmatrix} + \xi_j + \Delta\xi_{jt}, \text{and}$$

$$\mu_{ijt}(p_j,x_j,D_i,v_i;\Pi,\Sigma) = [-p_{jt},x_{j,1},x_{j,2},x_{j,3}]\left\{\begin{bmatrix} \Pi_{11} & \cdots & \Pi_{1,4} \\ \vdots & \ddots & \vdots \\ \Pi_{4,1} & \cdots & \Pi_{4,4} \end{bmatrix}\begin{bmatrix} D_{i,income} \\ D_{i,income^2} \\ D_{i,education} \\ D_{i,hhsize} \end{bmatrix}\right\}$$

$$+[-p_{jt},x_{j,1},x_{j,2},x_{j,3}]\begin{bmatrix} \Sigma_{11} & \cdots & \Sigma_{1,4} \\ \vdots & \ddots & \vdots \\ \Sigma_{4,1} & \cdots & \Sigma_{4,4} \end{bmatrix}\begin{bmatrix} v_{i,1} \\ v_{i,2} \\ v_{i,3} \\ v_{i,4} \end{bmatrix}$$

# References

[1] Berry, S, J.Levinsohn, and A.Pakes. 1995. "Automobile Prices in Market Equilibrium." *Econometrica* 63(4):841-890

[2] Berry, S, and A.Pakes. 2007."The Pure Characteristics Demand Model." *International Economic Review* 48(4):1193-1225

[3] Corinne, A., J.Balagtas, C.Mayen, and C.Greene. 2007. "Marketing Organic Milk in the United States: Findings from Agricultural Resource Management Survey of 2005." *Presentation at AAEA Annual Meeting.*

[4] Knittel, C. R. and K.Metaxoglou. 2008. "Estimation of Random Coefficient Demand Models: Challenges, Difficulties and Warnings." Working Paper, NBER No. W14080.

[5] Nevo, A. 1998. "A Research Assistant's Guide to Random Coefficients Discrete Choice Models of Demand." Technical Wroking Paper 221, NBER

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

[7] Nevo, A., and C.Wolfram. 2002. "Why Do Manufacturers Issue Coupons? An Empirical Analysis of Breakfast Cereals" *The RAND Journal of Economics*33(2):319-339

[8] Prasad,A., A.Strijnev, and Q.Zhang. 2008. "What can grocery basket data tell us about health consciousness?" *Journal of Research in Marketing.* 25:301-209

[9] Villas-Boas, J.M., and Y.Zhao. 2005. "Retailer, Manufacturers, and Individual Consumers: Modeling the Supply Side in the Ketchup Marketplace." *Journal of Marketing Research* 42(1):83-95

[10] Wathieu, L., and M. Bertini.2007. "Price as a Stimulus to Think: The Case for Willful Overpricing" *Marketing Science* 26(1):118–129

[11] Xiao, W. 2008. "The Competitive and Welfare Effects of New Product Introduction: The Case of Crystal Pepsi." Working Paper, Dept.of Economics., Johns Hopkins University.

Table 1: Volume Market Shares

|  | 2001 | 2002 | 2003 | 2004 | 2005 | 2006 |
|---|---|---|---|---|---|---|
| Heinz | 59.56% | 58.04% | 57.80% | 57.31% | 59.91% | 58.40% |
| Hunt's | 19.53% | 20.86% | 21.65% | 20.96% | 20.63% | 21.94% |
| Private Label | 15.75% | 16.19% | 15.09% | 16.86% | 14.99% | 15.61% |
| Del Monte | 4.62% | 4.18% | 4.87% | 4.13% | 3.76% | 3.37% |

source: IRI data base

Table 2: Statistical summary of Price and market share

|  | Mean | Median | Std | Min | Max |
|---|---|---|---|---|---|
| Price($/ounce) | 0.0748 | 0.0637 | 0.0396 | 0.0022 | 0.4297 |
| Potential Market Share(%) | 1.7606 | 0.5454 | 3.5023 | 0.0002 | 77.3240 |

source: IRI data base

Table 3: Logit Results

|  | OLS | | | 2SLS | | | | |
|---|---|---|---|---|---|---|---|---|
| Variable | i | ii | iii | iv | v | vi | vii | viii |
| Price | -16.59 | -6.45 | -6.448 | -19.99 | -2.89 | -3.01 | -2.55 | -13.22 |
| t-stat | -95.12 | -25.77 | -25.74 | -98.20 | -5.52 | -5.74 | -4.88 | -22.36 |
| Intercept | -1.74 | -4.68 | -7.73 | -1.47 | -8.56 | -11.42 | -11.46 | -7.94 |
| t-stat | -32.64 | -87.52 | -9.74 | -27.01 | -24.57 | -13.21 | -13.25 | -23.45 |
| Healthy | -1.68 | ___ | ___ | -1.45 | ___ | ___ | ___ | ___ |
| t-stat | -81.57 | ___ | ___ | -67.07 | ___ | ___ | ___ | ___ |
| Size | -0.26 | ___ | ___ | -0.41 | ___ | ___ | ___ | ___ |
| t-stat | -9.70 | ___ | ___ | -15.37 | ___ | ___ | ___ | ___ |
| log(mean(income)) | ___ | ___ | 0.17 | ___ | ___ | 0.16 | 0.16 | ___ |
| t-stat | ___ | ___ | 9.23 | ___ | ___ | 8.96 | 8.92 | ___ |
| log(mean(educ)) | ___ | ___ | 0.39 | ___ | ___ | 0.35 | 0.36 | ___ |
| t-stat | ___ | ___ | 1.71 | ___ | ___ | 1.55 | 1.57 | ___ |
| mean(hhsize) | ___ | ___ | -0.05 | ___ | ___ | -0.08 | -0.08 | ___ |
| t-stat | ___ | ___ | -1.71 | ___ | ___ | -2.44 | -2.54 | ___ |
| $R^2$ or 1st stage $R^2$ | 0.30 | 0.61 | 0.62 | 0.85 | 0.88 | 0.88 | 0.88 | 0.89 |
| 1st stage F test | ___ | ___ | ___ | 688.22 | 5001.15 | 4924.02 | 4804.79 | 4401.86 |
| Instruments | ___ | ___ | ___ | UPC dummies | prices | prices | prices, costs | prices, costs |

*Table 4: Brands and Number of UPC Used in Estimation*

| Regular | # of UPC | Organic | # of UPC | Reduced /No Sugar | # of UPC | No Salt | # of UPC |
|---|---|---|---|---|---|---|---|
| Heinz | 43 | Annie's Natural REG REGSL | 1 | Private Label | 4 | Heinz TMT NST | 1 |
| Hunt's | 17 | Heinz REG | 2 | Franks TMT REGS | 1 | Hunt's REG NSTA | 1 |
| Del Mont | 4 | Seed of Change TMT REGSL | 2 | Estee REG REGS | 1 | | |
| Private Label | 21 | Trees of Life REG REGSL | 2 | Heinz TMT REGSL | 1 | | |
| | | Muir Glen TMT REGSL | 1 | Heinz One Carb | 1 | | |
| | | | | Red Gold | 1 | | |
| | | | | Walden Farms | 1 | | |
| Total | 85 | Total | 8 | Total | 10 | Total | 2 |

Table 5: Statistics Summary for Demographics and Product Characteristics

| | Mean | Median | Std | Min | Max |
|---|---|---|---|---|---|
| Individual Income | 25,889.95 | 19,401.50 | 26,657.28 | 0.25 | 498,687 |
| Education Attainment | 39.81 | 40(College But No Degree) | 2.93 | 31(Less than 1st Grade) | 46(PHD,EDD) |
| Household size | 3.2 | 3 | 1.55 | 1 | 12 |
| Healthy | 0.13 | 0 | 0.33 | 0 | 1 |
| Size | 34.33 | 24 | 22.56 | 12 | 384 |