

Analysis of Household Buying Behavior: Carbonated Soft Drinks and Other Beverages

Günter J. Hitsch

1/15/2017

We will perform an analysis of category buying and consumption behavior of beverages (not including alcohol or dairy), with a particular emphasis on sodas (carbonated soft drinks/CSD's). Sodas have been the subject of an intense debate linked to public health concerns, and one or two cent-per-ounce taxes have recently been passed in San Francisco, Cook County (IL), Philadelphia, Boulder (CO), and other places. Changes in buying behavior due to health concerns present challenges to the beverage manufacturers, but also opportunities if consumers are shifting their consumption to healthier substitutes.

Data

The Nielsen Homescan panel data set is an ideal data source to study broad trends in consumption behavior. This data set is available for research and teaching purposes from the Kilts Center for Marketing at Chicago Booth.

Currently, the Homescan panel at the Kilts Center contains data from 2004 to 2014. In many years we have information on the buying behavior of more than sixty thousand households. The corresponding full data set is large. Hence, to avoid memory and computing time issues, I extracted the beverage data that we will use for the analysis. Consult the `Buying-Behavior-Data-Preparation.Rmd` file to see the steps taken in the extraction and cleaning process. Once you load the data, you can verify that even the beverage data only contains more than 40 million observations (use `nrow`, `ncol`, or `dim` to see the size of a `data.table` or `data.frame`).

The key data for the analysis are contained in a **purchase file** and a **product file**. Load the data:

```
library(bit64)
library(data.table)

data_folder    = "./Data"
purchases_file = "purchases_beverages.RData"
products_file  = "products_beverages.RData"

load(paste0(data_folder, "/", purchases_file))
load(paste0(data_folder, "/", products_file))
```

Note that the data are in a sub-folder of *this* R Markdown document called **Data**, i.e. they are *assumed* to be in that folder! The `paste0` command merges strings. The `bit64` package is used because the product UPC numbers are 64-bit (long) integers, a data type that is not part of base R.

Using a subsample

Important trick: To develop and test your code, you can work with a subsample of your original data. For example, to draw a random sample of 4 million observations without replacement, use:

```
purchases_sub = purchases[sample(.N, 4000000)]
```

For details, consult `?sample`, and note that `.N` is the number of rows in a `data.table`—a variable that the `data.table` package automatically supplies.

Alternatively, it may be even better to obtain *all* purchase data for a random sample of households. First, we obtain a 25 percent sample of all household codes in the data:

```
N_households = length(unique(purchases$household_code))
N_subsample   = round(0.25*N_households)

household_code_sub = sample(unique(purchases$household_code), N_subsample)
```

Then extract all data for the chosen household keys. As we already discussed in the data.table *Keys and Merging* overview, the first method is more readable yet somewhat slower, the second method is faster but also more confusing to the novice. The second method also does not keep it's key, so you have to key the `purchases_sub_hh_a` data.table later.

```
purchases_sub_hh   = purchases[household_code %in% household_code_sub]
purchases_sub_hh_a = purchases[.(household_code_sub)]
```

For this analysis, you may just use the original data (40 million observations is quite easily manageable using data.table on a fast computer), but it's nonetheless an important trick to remember! Also, it typically makes sense to create the subsample in an initial step and save it to a file.

To conserve memory we remove (delete) all three subsamples:

```
rm(purchases_sub)
rm(list = c("purchases_sub_hh", "purchases_sub_hh_a"))
```

Variable description

The **purchases** data.table contains information on the products bought by the households.

- **Inspect the data.**

Households are identified based on a unique `household_code`. For each shopping trip we know the `purchase_date` and the `retailer_code` (for confidentiality, the Kilts Center data do not include the exact name of the retailer).

For each product we have information on the `total_price_paid` and the `quantity` purchased. A deal flag (0/1) and a `coupon_value` is also provided. The `total_price_paid` applies to *all* units purchased. Hence, if `total_price_paid` = 7.98 and `quantity` = 2, then the *per-unit price* is $7.98/2 = 3.99$ dollars. Furthermore, if the `coupon_value` is positive, then the total dollar amount that the household spent is `total_price_paid` - `coupon_value`. The *per-unit cost* to the household is then even lower. For example, if `coupon_value` = 3 in the example above, then the per-unit cost to the household is $(7.98 - 3)/2 = 2.49$ dollars.

I recommend to use the per-unit cost if the objective is to measure household dollar spending per unit purchased. If the objective is to measure the shelf-price of the product, use the per-unit price instead.

Important data notes:

1. Products in the Nielsen Homescan and RMS scanner data are identified by a unique combination of `upc` and `upc_ver_uc`. Why not just the UPC code? — Because UPC's can change over time, for example if a UPC is assigned to a different product. The UPC version code captures such a change. From now on, whenever we refer to a *product*, we mean a `upc/upc_ver_uc` combination. To identify a unique product in data.table, use a `by = .(upc, upc_ver_uc)` statement.
2. The `panel_year` variable in **purchases** is intended *only* to link households to the corresponding household attribute data (income, age, etc.), which are updated yearly. At the beginning of a `panel_year` it need not exactly correspond to the calendar year.

The **products** data.table contains product information for each `upc/upc_ver_uc` combination.

- **Inspect the data.table.**

Note that products are organized into departments (e.g. DRY GROCERY), product groups (e.g. CARBONATED BEVERAGES), and product modules (e.g. SOFT DRINKS - CARBONATED), with corresponding codes and descriptions. Use `unique` or `table` to print all values for the variables. Or create a table that lists all the product module codes, product module descriptions, and group descriptions in the data:

```
module_DT = products[, head(.SD, 1), by = product_module_code,
                          .SDcols = c("product_module_descr", "product_group_descr")]
module_DT[order(product_group_descr)]
```

We also have brand codes and descriptions, such as PEPSI R (Pepsi regular). You will often see the brand description CTL BR, which stands for *control brand*, i.e. private label brand. Brands are identified using either the `brand_code_uc` or `brand_descr` variables, and are sold as different products (`upc/upc_ver_uc`) that differ along size, form (e.g. bottles vs. cans), or flavor.

`multi` indicates the number of units in a multi-pack (a multi-pack is a pack size such that `multi > 1`). More on the amount and unit variables below.

For many more details on the data, consult the *Consumer Panel Dataset Manual* (on Canvas).

Prepare the data for the analysis

We will calculate yearly summary statistics of customer buying behavior. To calculate the year corresponding to a purchase date, use `year()` in the data.table `IDateTime` class (see `?IDateTime` to learn about other conversion functions).

```
purchases[, year := year(purchase_date)]
```

Note that we create this year-variable because the `panel_year` variable in `purchases` does not exactly correspond to the calendar year, as we already discussed above. You can verify that there is a tiny percentage of observations for the 2003 calendar year, that “slipped” into the data set because the purchases are recorded for households in the 2004 `panel_year`.

- **Remove the 2003 observations.**

Define categories

In the analysis we want to distinguish between carbonated soft drinks (CSD’s), diet (low-calorie) CSD’s, bottled water, and a category including all other beverages (juices, ...). Hence, we create a new `category` variable that assigns each beverage purchase to one of these four categories.

- **First, create a default category variable with a name such as “Other”. Then find the product module codes for the three relevant categories and assign a corresponding name to category.**
- **Document the number of observations in each category.**

Now merge the category variable with the purchase data.

```
purchases = merge(purchases, products[, .(upc, upc_ver_uc, category)])
```

Note that `by = .(upc, upc_ver_uc, ...)` provides the product-level link between the products and the purchase table.

Volume in equivalent units

To measure volume in *equivalent units*, we need the product-level information on the *units of measurement* of product volume (`size1_units`), such as ounces, and the corresponding volume in a pack

size (size1_amount1). This information needs to be merged with the purchase data. We also merge with multi', which indicates multi-pack sizes.

- **Perform this merge.**

For beverages, product volume is typically measured in ounces (OZ), less frequently in quarts (QT), and only rarely in counts (CT).

- **Document the number of observations by unit of measurement. Let's ignore counts, and remove all corresponding data from the purchases data.table. Then convert the quantity of units purchased into a common volume measure in gallons, using the size1_amount variable. Also incorporate the multi variable into the volume calculation—multi accounts for multi-packs.**

Number of households in the data

To calculate the number of households in the data by year, use:

```
purchases[, no_households := length(unique(household_code)), by = year]
```

- **Create and show a table with the number of households by year.**

Note the expansion in the number of Homescan panelists in 2007!

Category-level analysis

Now we are ready to analyse the evolution of purchases and consumption in the four product categories. We want to calculate total and per capita consumption metrics. First, we create the total dollar spend and the total purchase volume for each category/year combination. We use the data.table approach for aggregation:

```
purchases_category = purchases[, .(spend           = sum(total_price_paid - coupon_value),  
  purchase_volume = sum(volume),  
  no_households   = head(no_households, 1)),  
  keyby = .(category, year)]
```

- **Calculate per capita spend and purchase volume (in gallons) for each category separately. Then graph the evolution of the yearly per capita purchase volume for all four categories.**

Note: Instead of creating graphs for each of the four categories you can use a facet_wrap layer provided by ggplot2.

- **Express the purchase/consumption data as multiples of the 2004 values, such that per capita volume takes the value of 1.0 in all categories in 2004. Such a normalization allows us to compare the consumption series in each category directly in percentage terms. Then show the graphs of consumption (normalized to its 2004 value), and discuss the results.**

Brand-level analysis

Now let's investigate the evolution of consumption for some of the key brands in the soda and bottled water categories.

- **First, merge the brand identifier brand_descr with the purchase data.**

Then we rank brands by total dollar spend in each category separately. We can assign ranks either using the rank function in base R, or using frankv (or frank) in the data.table package. Simple usage: To rank a vector x in ascending order, use frankv(x). To rank in descending order, use frankv(x, order = -1). See ?frank for more options.

First calculate total dollar spend by each category/brand combination, then assign the rank according to total spend:

```
brand_summary = purchases[, .(spend = sum(total_price_paid - coupon_value)),
                             by = .(category, brand_descr)]
brand_summary[, rank := frankv(spend, order = -1), by = category]
```

- Merge the brand ranks in the brand_summary table with the purchases information. Aggregate to the brand level, as we did before at the category level. Then calculate per capita and normalize per capita consumption. Plot the evolution of brand volume for the top 4 brands, separately for the CSD, diet CSD, and bottled water categories.

For bottled water you may add the `scales = "free_y"` option in `facet_wrap`:

```
facet_wrap(..., scales = "free_y")
```

By default, the y-axes in a facet wrap are identical across all panels, but `free_y` let's ggplot2 choose different axes for each panel.

- Discuss your results.