

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

Giovanni Compiani

Contents

1	Overview	2
2	Data	2
3	Variable description	3
4	Prepare the data for the analysis	5
4.1	Define categories	5
4.2	Volume in equivalent units	6
4.3	Number of households in the data	7
5	Category-level analysis	8
6	Brand-level analysis	11
7	Discussion	12
8	Appendix: Using a random subsample of the data	13

1 Overview

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and one or two cent-per-ounce taxes have recently been passed in San Francisco, Cook County (although the tax was repealed in 2017), Philadelphia, Boulder, and other cities. 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.

2 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.

We will use Homescan panel 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. I also took a *25 percent random subsample* of all the original observations. Once you load the data, you can verify that even this subsample of the beverage data contains more than 10 million observations (use `nrow`, `ncol`, or `dim` to see the size of a `data.table` or `data frame`).

Although not necessary for this assignment, it is good to know how to create random subsamples. For more information, please consult the Appendix below.

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

```
library(bit64)
library(data.table)
library(ggplot2)

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 string `data_folder` contains the location of the data, i.e. the data 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.

3 Variable description

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

- Inspect the data

```
# Inspect structure of the data
str(purchases)

# View first few rows of the data
head(purchases)
```

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. We will work with the household attribute data in one of the next assignments.

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

- Inspect the data.table

```
# Inspect structure of the data
str(products)

# View first few rows of the data
head(products)
```

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).

4 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). Note that there are other methods for time aggregation that we will study later in the course.

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

Key: <household_code, purchase_date, upc, upc_ver_uc>

	upc	upc_ver_uc	household_code	retailer_code	purchase_date
	<i64>	<int>	<int>	<int>	<Date>
1:	3200000046	1	2000021	6600	2004-02-07
2:	83281100322	1	2000021	6600	2004-05-12
3:	5260309050	1	2000021	6999	2004-06-05
4:	80276311327	1	2000021	6999	2004-07-31
5:	8489300900	1	2000021	5852	2004-08-14
6:	8489300900	1	2000021	5852	2004-09-19

	total_price_paid	quantity	coupon_value	deal_flag_uc	panel_year	year
	<num>	<int>	<num>	<int>	<int>	<int>
1:	0.38	1	0	0	2004	2004
2:	0.66	1	0	0	2004	2004
3:	0.99	1	0	0	2004	2004
4:	0.99	1	0	0	2004	2004
5:	1.98	2	0	0	2004	2004
6:	2.97	3	0	0	2004	2004

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

```
# Remove observations where the year is 2003
purchases <- purchases[year != 2003]

# Check if the 2003 observations were successfully removed
summary(purchases$year)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
2004	2007	2009	2009	2012	2014

4.1 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 allows us classify the beverage purchase observations.

- First, in the `products` table, 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, i.e. the number of products that belong to each of the categories.

```
# Assuming 'products' is your data.table
```

```
# Step 1: Create a default category called 'Other'
products[, category := "Other"]

# Step 2: Assign specific categories based on 'product_module_descr'
products[product_module_code == "1484", category := "Carbonated Soft Drinks"]
products[product_module_code == "1553", category := "Diet Soft Drinks"]
products[product_module_code == "1487", category := "Bottled Water"]

# Step 3: Check the distribution of categories
# table(products$category)
products[, .N, by = category]
```

```
      category      N
      <char> <int>
1:      Other 58754
2: Carbonated Soft Drinks 24126
3:      Diet Soft Drinks 10267
4:      Bottled Water 15921
```

Now merge the category variable with the purchase data.

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

To understand why this statement works, inspect the keys of the two tables (e.g. `key(purchases)`). The keys contain two common columns, `upc` and `upc_ver_uc`, which are used as the foreign key that provides the product-level link between the products and purchase tables.

4.2 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**

```
purchase1 = merge(purchases, products[, .(upc, upc_ver_uc, size1_units, size1_amount, multi)])
```

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.****

Hint: For example, for observations where `size1_units` is OZ, you can calculate the volume:

```
volume := (1/128)*size1_amount*multi*quantity

# ignore counts and remove all corresponding data from the purchases data. table
purchase1 = purchase1[size1_units != "CT"]

# Inspect the data
unique_values = unique(purchase1$size1_units)
print(unique_values)
```

```
[1] "OZ" "QT"
```

```
# convert the `quantity` of units purchased into a common `volume` measure in gallons, using the `size`
purchase1[, volume := ifelse(size1_units == "OZ",
                             (1/128) * size1_amount * multi * quantity,
                             ifelse(size1_units == "QT",
                                     (1/4) * size1_amount * multi * quantity,
                                     NA))] # NA for other units
```

4.3 Number of households in the data

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

```
purchase1[, 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!

```
# Calculate the number of unique households by year
households_by_year = purchase1[, .(no_households = length(unique(household_code))), by = year]

# Display the result as a table
print(households_by_year)
```

	year	no_households
	<int>	<int>
1:	2012	15395
2:	2014	15284
3:	2009	16077
4:	2010	16322
5:	2007	16330
6:	2013	15764
7:	2011	15366
8:	2005	9700
9:	2004	10566
10:	2008	16186
11:	2006	10565

5 Category-level analysis

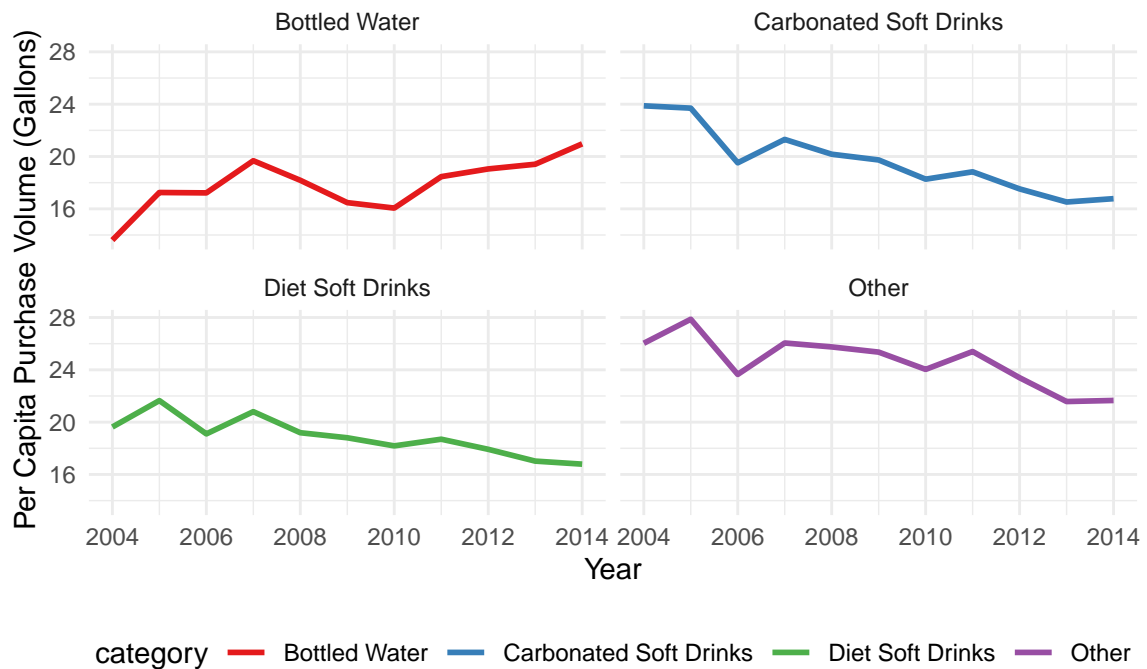
Now we are ready to analyse the evolution of purchases and consumption in the four product categories. We want to calculate total and per capita (more precisely: per household) consumption metrics. First, we create the total dollar spend and the total purchase volume for each category/year combination.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

```
purchases_category = purchase1[,  
  .(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.

```
# Calculate per capita spend and purchase volume  
purchases_category[, `:=`(  
  per_capita_spend = spend / no_households,  
  per_capita_volume = purchase_volume / no_households  
)]  
  
# Plot the evolution of per capita purchase volume using ggplot2  
  
ggplot(purchases_category, aes(x = year, y = per_capita_volume, color = category)) +  
  geom_line(size = 1) +  
  labs(title = "Evolution of Per Capita Purchase Volume by Category",  
    x = "Year",  
    y = "Per Capita Purchase Volume (Gallons)") +  
  theme_minimal() +  
  scale_color_brewer(palette = "Set1") +  
  theme(legend.position = "bottom") +  
  facet_wrap(~ category) # Separate graphs by category
```


Evolution of Per Capita Purchase Volume by Category



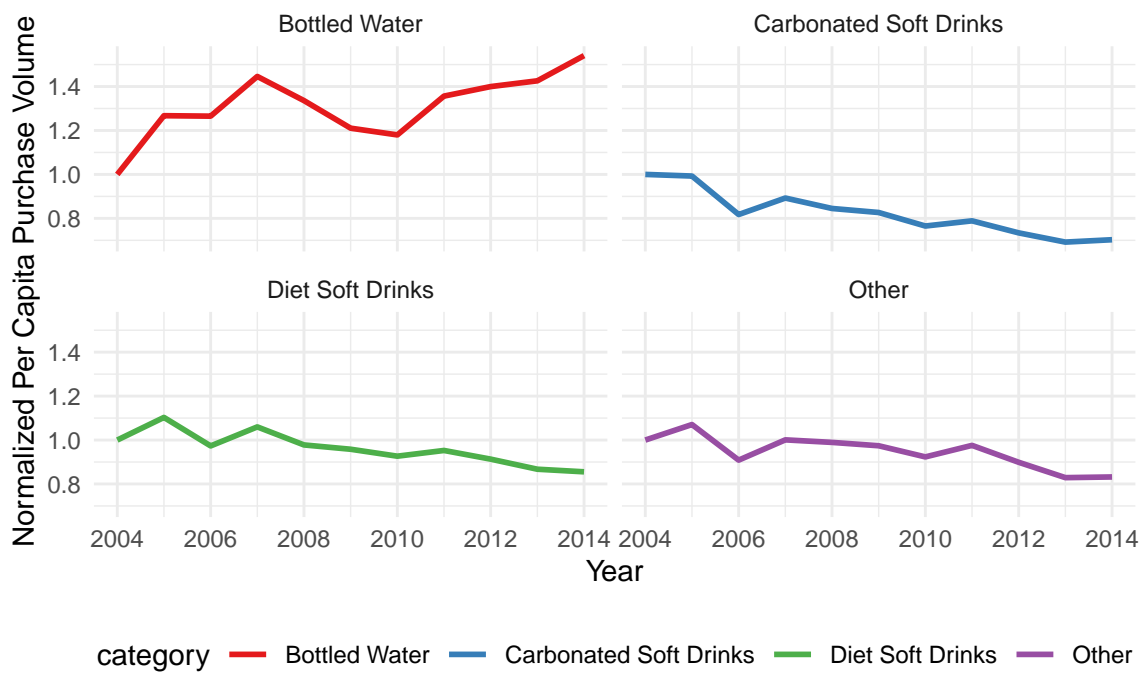
Note: Instead of creating graphs for each of the four categories you can use a `facet_wrap` layer provided by `ggplot2`.

- 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.

```
# Normalize per capita purchase volume to 2004 values
purchases_category[, normalized_per_capita_volume := per_capita_volume / per_capita_volume[year == 2004]]

# Plot the normalized per capita purchase volume
ggplot(purchases_category, aes(x = year, y = normalized_per_capita_volume, color = category)) +
  geom_line(size = 1) +
  labs(title = "Evolution of Normalized Per Capita Purchase Volume by Category (2004 = 1.0)",
       x = "Year",
       y = "Normalized Per Capita Purchase Volume") +
  theme_minimal() +
  scale_color_brewer(palette = "Set1") +
  theme(legend.position = "bottom") +
  facet_wrap(~ category) # Separate graphs by category
```

Evolution of Normalized Per Capita Purchase Volume by Category (2)



6 Brand-level analysis

Now we 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

```
# Merge the brand identifier `brand_descr` with the purchase data
purchase2 = merge(purchase1, products[, .(upc, upc_ver_uc, brand_descr)])
```

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 = purchase2[, .(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 spending and volume, and normalize the per capita variables to 1.0 in 2004, as before at the category level. Plot the evolution of brand volume for the top four brands, separately for the CSD, diet CSD, and bottled water categories.

```
# Merge the brand ranks in the `brand_summary` table with the `purchases` information.
purchase2 = merge(purchase2, brand_summary[, .(category, brand_descr, rank)], by = c("category", "brand_

# Aggregate to the brand level
purchases_brand = purchase2[,
  .(spend      = sum(total_price_paid - coupon_value),
    purchase_volume = sum(volume),
    no_households  = head(no_households, 1),
    rank= min(rank)),
  keyby = .(category, brand_descr, year)]

# Calculate per capita spending and volume
purchases_brand[, `:=`(
  per_capita_spend = spend / no_households,
  per_capita_volume = purchase_volume / no_households
)]

# Normalize per capita purchase volume to 2004 values at the category level
purchases_brand[, `:=`(
  normalized_per_capita_spend = per_capita_spend / per_capita_spend[year == 2004],
  normalized_per_capita_volume = per_capita_volume / per_capita_volume[year == 2004]
), by = .(category, brand_descr)]
```

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

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

```
top_4_brands = purchases_brand[rank <= 4]
```

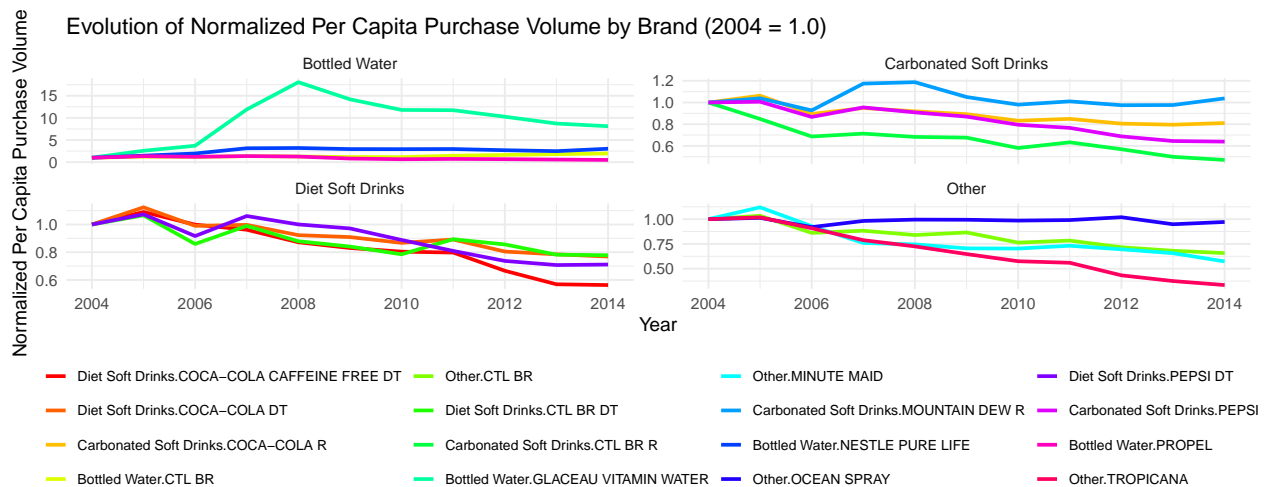
```
# Generate 16 random colors
```

```

set.seed(123) # Set seed for reproducibility
random_colors <- rainbow(16)

# Use the random colors in your ggplot
ggplot(top_4_brands, aes(x = year, y = normalized_per_capita_volume, color = interaction(category, brand))) +
  geom_line(size = 1) +
  labs(title = "Evolution of Normalized Per Capita Purchase Volume by Brand (2004 = 1.0)",
       x = "Year",
       y = "Normalized Per Capita Purchase Volume") +
  theme_minimal() +
  scale_color_manual(values = random_colors) + # Use the random color palette
  theme(legend.position = "bottom",
       legend.key.size = unit(0.5, 'cm'),
       legend.text = element_text(size = 8)) +
  guides(color = guide_legend(title = NULL)) +
  facet_wrap(~ category, scales = "free_y") # Separate graphs by category

```



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

7 Discussion

Provide a brief discussion of the marketing implications of your findings.

Bottled water consumption showed a strong upward trend from 2004 to 2014, with notable spikes in 2007 and 2012-2014. This trend suggests a shift towards healthier options, presenting an opportunity for brands to emphasize health benefits and hydration, particularly for active consumers.

Conversely, both carbonated and diet soft drinks declined. Regular soft drinks started decreasing after 2007, while diet drinks peaked around 2006 before declining. This likely reflects growing concerns over sugar and artificial ingredients.

Premium brands like Nestle Pure Life and Vitamin Water performed well, especially post-2007, indicating consumer preference for differentiated products.

The beverage market from 2007 to 2014 clearly shifted towards healthier options, with bottled water emerging as a key beneficiary. To capitalize on these trends, soft drink brands need to innovate with healthier alternatives or smaller portions. Premium water brands should continue product innovation, such as adding functional benefits. All brands must adapt to this health-conscious market landscape to remain competitive.

8 Appendix: Using a random subsample of the data

Although not necessary for the analysis in this assignment, especially as I already created a random subsample of the original data, it is useful to know how to create such a random subsample.

For example, to draw a random sample of 2 million observations without replacement, use:

```
purchases_sub = purchases[sample(.N, 2000000)]
```

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 its 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)]
```