

# Quantium Virtual Internship - Retail Strategy and Analytics - Task

1

## Solution for Task 1

### Load required libraries and datasets

```
#### Example code to install packages
#install.packages("data.table")
#### Load required libraries
library(data.table)
library(ggplot2)
library(ggmosaic)
library(readr)

library(dplyr)
library(tidyr)
library(tidytext)
library(stringr)
library(ggforce)
library(purrr)
library(grid)
library(arules)
library(RColorBrewer)

#### Point the filePath to where you have downloaded the datasets to and
#### assign the data files to data.tables
# over to you! fill in the path to your working directory.
filePath <- "~/Desktop/Forage/Quantium-Data_Analytics/"
transactionData <- fread(paste0(filePath, "QVI_transaction_data.csv"))
customerData <- fread(paste0(filePath, "QVI_purchase_behaviour.csv"))
```

### Exploratory data analysis

The first step in any analysis is to first understand the data. Let's take a look at each of the datasets provided.

#### Examining transaction data

We can use `str()` to look at the format of each column and see a sample of the data. As we have read in the dataset as a `data.table` object, we can also run `transactionData` in the console to see a sample of the data or use `head(transactionData)` to look at the first 10 rows.

Let's check if columns we would expect to be numeric are in numeric form and date columns are in date format.

```
#### Examine transaction data
# Over to you! Examine the data using one or more of the methods described above.
str(transactionData)
```

```
## Classes 'data.table' and 'data.frame': 264836 obs. of 8 variables:
## $ DATE : int 43390 43599 43605 43329 43330 43604 43601 43601 43332 43330 ...
## $ STORE_NBR : int 1 1 1 2 2 4 4 4 5 7 ...
## $ LYLTY_CARD_NBR: int 1000 1307 1343 2373 2426 4074 4149 4196 5026 7150 ...
## $ TXN_ID : int 1 348 383 974 1038 2982 3333 3539 4525 6900 ...
## $ PROD_NBR : int 5 66 61 69 108 57 16 24 42 52 ...
## $ PROD_NAME : chr "Natural Chip Compny SeaSalt175g" "CCs Nacho Cheese 175g"
"Smiths Crinkle Cut Chips Chicken 170g" "Smiths Chip Thinly S/Cream&Onion 175g"
...
## $ PROD_QTY : int 2 3 2 5 3 1 1 1 1 2 ...
## $ TOT_SALES : num 6 6.3 2.9 15 13.8 5.1 5.7 3.6 3.9 7.2 ...
## - attr(*, ".internal.selfref")=<externalptr>
```

```
str(customerData)
```

```
## Classes 'data.table' and 'data.frame': 72637 obs. of 3 variables:
## $ LYLTY_CARD_NBR : int 1000 1002 1003 1004 1005 1007 1009 1010 1011 1012 ...
## $ LIFESTAGE : chr "YOUNG SINGLES/COUPLES" "YOUNG SINGLES/COUPLES" "YOUNG
FAMILIES" "OLDER SINGLES/COUPLES" ...
## $ PREMIUM_CUSTOMER: chr "Premium" "Mainstream" "Budget" "Mainstream" ...
## - attr(*, ".internal.selfref")=<externalptr>
```

```
head(transactionData)
```

```
##      DATE STORE_NBR LYLTY_CARD_NBR TXN_ID PROD_NBR
## 1: 43390          1          1000      1         5
## 2: 43599          1          1307     348        66
## 3: 43605          1          1343     383        61
## 4: 43329          2          2373     974        69
## 5: 43330          2          2426    1038       108
## 6: 43604          4          4074    2982        57
##
##      PROD_NAME PROD_QTY TOT_SALES
## 1:  Natural Chip      Compny SeaSalt175g      2      6.0
## 2:           CCs Nacho Cheese      175g      3      6.3
## 3:  Smiths Crinkle Cut Chips Chicken 170g      2      2.9
## 4:  Smiths Chip Thinly S/Cream&Onion 175g      5     15.0
## 5: Kettle Tortilla ChpsHny&Jlpno Chili 150g      3     13.8
## 6: Old El Paso Salsa  Dip Tomato Mild 300g      1      5.1
```

```
head(customerData)
```

```
##      LYLTY_CARD_NBR      LIFESTAGE PREMIUM_CUSTOMER
## 1:          1000 YOUNG SINGLES/COUPLES      Premium
## 2:          1002 YOUNG SINGLES/COUPLES      Mainstream
## 3:          1003      YOUNG FAMILIES      Budget
## 4:          1004 OLDER SINGLES/COUPLES      Mainstream
## 5:          1005 MIDAGE SINGLES/COUPLES      Mainstream
## 6:          1007 YOUNG SINGLES/COUPLES      Budget
```

We can see that the date column is in an integer format. Let's change this to a date format.

```
#### Convert DATE column to a date format
#### A quick search online tells us that CSV and Excel integer dates begin on 30 Dec 1899
```

```
transactionData$DATE <- as.Date(transactionData$DATE, origin = "1899-12-30")
```

We should check that we are looking at the right products by examining PROD\_NAME.

```
#### Examine PROD_NAME
# Over to you! Generate a summary of the PROD_NAME column.
summary(transactionData$PROD_NAME)
```

```
##      Length      Class      Mode
##      264836 character character
```

Looks like we are definitely looking at potato chips but how can we check that these are all chips? We can do some basic text analysis by summarising the individual words in the product name.

```
#### Examine the words in PROD_NAME to see if there are any incorrect entries
#### such as products that are not chips
productWords <- data.table(unlist(strsplit(unique(transactionData[, PROD_NAME]), " ")))
setnames(productWords, 'words')
```

As we are only interested in words that will tell us if the product is chips or not, let's remove all words with digits and special characters such as '&' from our set of product words. We can do this using `grepl()`.

```
# Over to you! Remove digits, and special characters, and then sort the distinct words by frequency of

#### Removing digits
digits <- !grepl('[:digit:]', productWords$words)
productWords <- productWords[digits,]
#### Removing special characters
special <- !grepl('[:punct:]', productWords$words)
productWords <- productWords[special,]
#### Let's look at the most common words by counting the number of times a word appears and
#### sorting them by this frequency in order of highest to lowest frequency

# product_wordcounts <- productWords %>%
#   count(words, sort = TRUE)
product_wordcounts <- productWords[, .N, words][order(desc(N))]
```

There are salsa products in the dataset but we are only interested in the chips category, so let's remove these.

```
#### Remove salsa products
transactionData[, SALSA := grepl("salsa", tolower(PROD_NAME))]
transactionData <- transactionData[SALSA == FALSE, ][, SALSA := NULL]
```

Next, we can use `summary()` to check summary statistics such as mean, min and max values for each feature to see if there are any obvious outliers in the data and if there are any nulls in any of the columns (NA's : number of nulls will appear in the output if there are any nulls).

```
#### Summarise the data to check for nulls and possible outliers
# Over to you!
summary(transactionData)
```

```
##      DATE      STORE_NBR  LYLTY_CARD_NBR      TXN_ID
##  Min.   :2018-07-01   Min.    : 1.0      Min.    : 1000   Min.    :    1
## 1st Qu.:2018-09-30   1st Qu.: 70.0     1st Qu.: 70015   1st Qu.: 67569
## Median :2018-12-30   Median :130.0     Median : 130367   Median : 135183
## Mean   :2018-12-30   Mean   :135.1     Mean   : 135531   Mean   : 135131
## 3rd Qu.:2019-03-31   3rd Qu.:203.0     3rd Qu.: 203084   3rd Qu.: 202654
## Max.   :2019-06-30   Max.    :272.0     Max.    :2373711   Max.    :2415841
```

##	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
##	Min. : 1.00	Length:246742	Min. : 1.000	Min. : 1.700
##	1st Qu.: 26.00	Class :character	1st Qu.: 2.000	1st Qu.: 5.800
##	Median : 53.00	Mode :character	Median : 2.000	Median : 7.400
##	Mean : 56.35		Mean : 1.908	Mean : 7.321
##	3rd Qu.: 87.00		3rd Qu.: 2.000	3rd Qu.: 8.800
##	Max. :114.00		Max. :200.000	Max. :650.000

There are no nulls in the columns but product quantity appears to have an outlier which we should investigate further. Let's investigate further the case where 200 packets of chips are bought in one transaction.

```
#### Filter the dataset to find the outlier
# Over to you! Use a filter to examine the transactions in question.
transactionData %>% filter(PROD_QTY == 200)
```

##	DATE	STORE_NBR	LYLTY_CARD_NBR	TXN_ID	PROD_NBR
##	1: 2018-08-19	226	226000	226201	4
##	2: 2019-05-20	226	226000	226210	4

##	PROD_NAME	PROD_QTY	TOT_SALES
##	1: Dorito Corn Chp Supreme 380g	200	650
##	2: Dorito Corn Chp Supreme 380g	200	650

There are two transactions where 200 packets of chips are bought in one transaction and both of these transactions were by the same customer.

```
#### Let's see if the customer has had other transactions
# Over to you! Use a filter to see what other transactions that customer made.
transactionData %>% filter(LYLY_CARD_NBR == 226000)
```

##	DATE	STORE_NBR	LYLTY_CARD_NBR	TXN_ID	PROD_NBR
##	1: 2018-08-19	226	226000	226201	4
##	2: 2019-05-20	226	226000	226210	4

##	PROD_NAME	PROD_QTY	TOT_SALES
##	1: Dorito Corn Chp Supreme 380g	200	650
##	2: Dorito Corn Chp Supreme 380g	200	650

It looks like this customer has only had the two transactions over the year and is not an ordinary retail customer. The customer might be buying chips for commercial purposes instead. We'll remove this loyalty card number from further analysis.

```
#### Filter out the customer based on the loyalty card number
# Over to you!
transactionData <- transactionData %>%
  filter(LYLY_CARD_NBR != 226000)
```

```
#### Re-examine transaction data
# Over to you!
summary(transactionData)
```

##	DATE	STORE_NBR	LYLTY_CARD_NBR	TXN_ID
##	Min. :2018-07-01	Min. : 1.0	Min. : 1000	Min. : 1
##	1st Qu.:2018-09-30	1st Qu.: 70.0	1st Qu.: 70015	1st Qu.: 67569
##	Median :2018-12-30	Median :130.0	Median : 130367	Median : 135182
##	Mean :2018-12-30	Mean :135.1	Mean : 135530	Mean : 135130
##	3rd Qu.:2019-03-31	3rd Qu.:203.0	3rd Qu.: 203083	3rd Qu.: 202652
##	Max. :2019-06-30	Max. :272.0	Max. :2373711	Max. :2415841

##	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
##	Min. : 1.00	Length:246740	Min. :1.000	Min. : 1.700

```
## 1st Qu.: 26.00    Class :character    1st Qu.:2.000    1st Qu.: 5.800
## Median : 53.00    Mode  :character    Median :2.000    Median : 7.400
## Mean   : 56.35                    Mean   :1.906    Mean   : 7.316
## 3rd Qu.: 87.00                    3rd Qu.:2.000    3rd Qu.: 8.800
## Max.   :114.00                    Max.    :5.000    Max.    :29.500
```

That's better. Now, let's look at the number of transaction lines over time to see if there are any obvious data issues such as missing data.

```
#### Count the number of transactions by date
# Over to you! Create a summary of transaction count by date.
transactionData[, .N, DATE][order(DATE)]
```

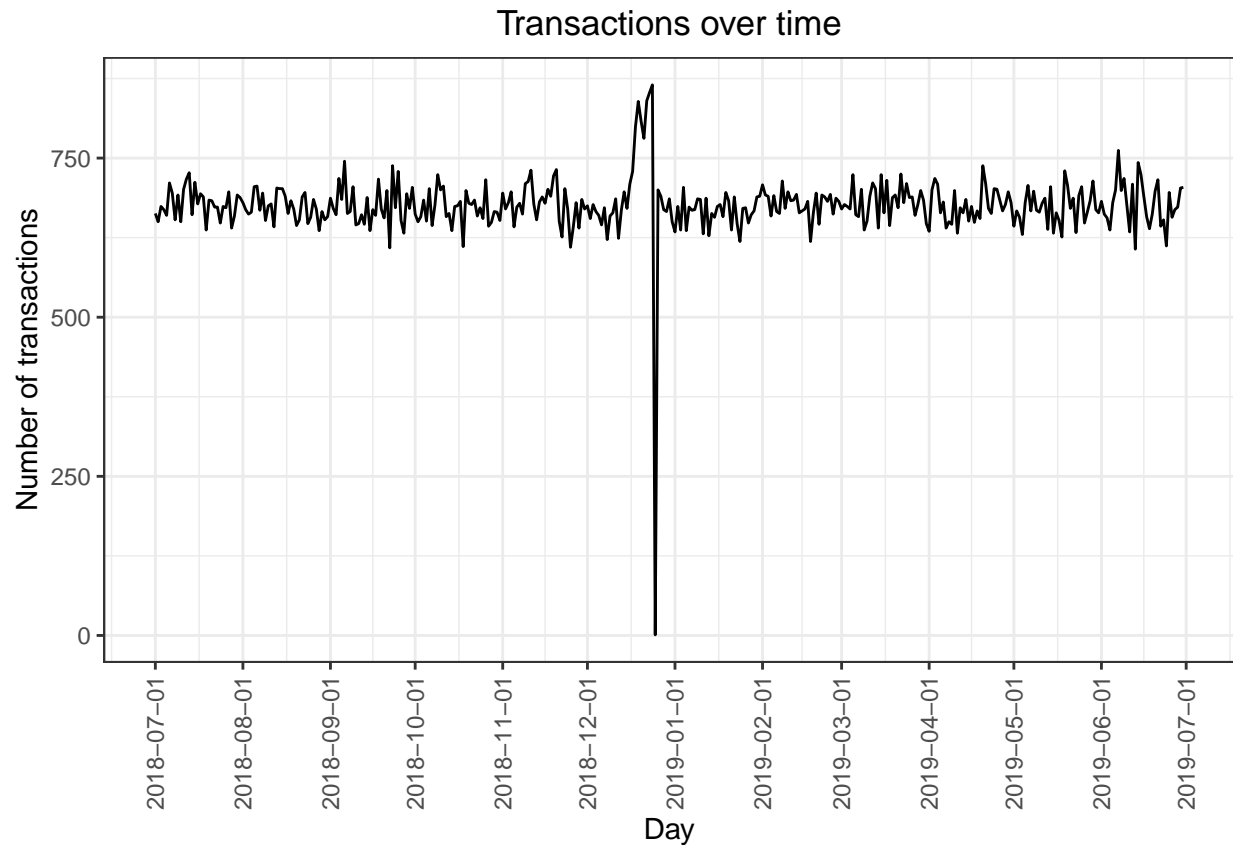
```
##          DATE      N
## 1: 2018-07-01 663
## 2: 2018-07-02 650
## 3: 2018-07-03 674
## 4: 2018-07-04 669
## 5: 2018-07-05 660
## ---
## 360: 2019-06-26 657
## 361: 2019-06-27 669
## 362: 2019-06-28 673
## 363: 2019-06-29 703
## 364: 2019-06-30 704
```

There's only 364 rows, meaning only 364 dates which indicates a missing date. Let's create a sequence of dates from 1 Jul 2018 to 30 Jun 2019 and use this to create a chart of number of transactions over time to find the missing date.

```
#### Create a sequence of dates and join this the count of transactions by date
# Over to you - create a column of dates that includes every day from 1 Jul 2018 to 30 Jun 2019, and join
all_dates <- data_frame(DATE=seq(as.Date('2018-07-01'), as.Date('2019-06-30'), by='days'))
```

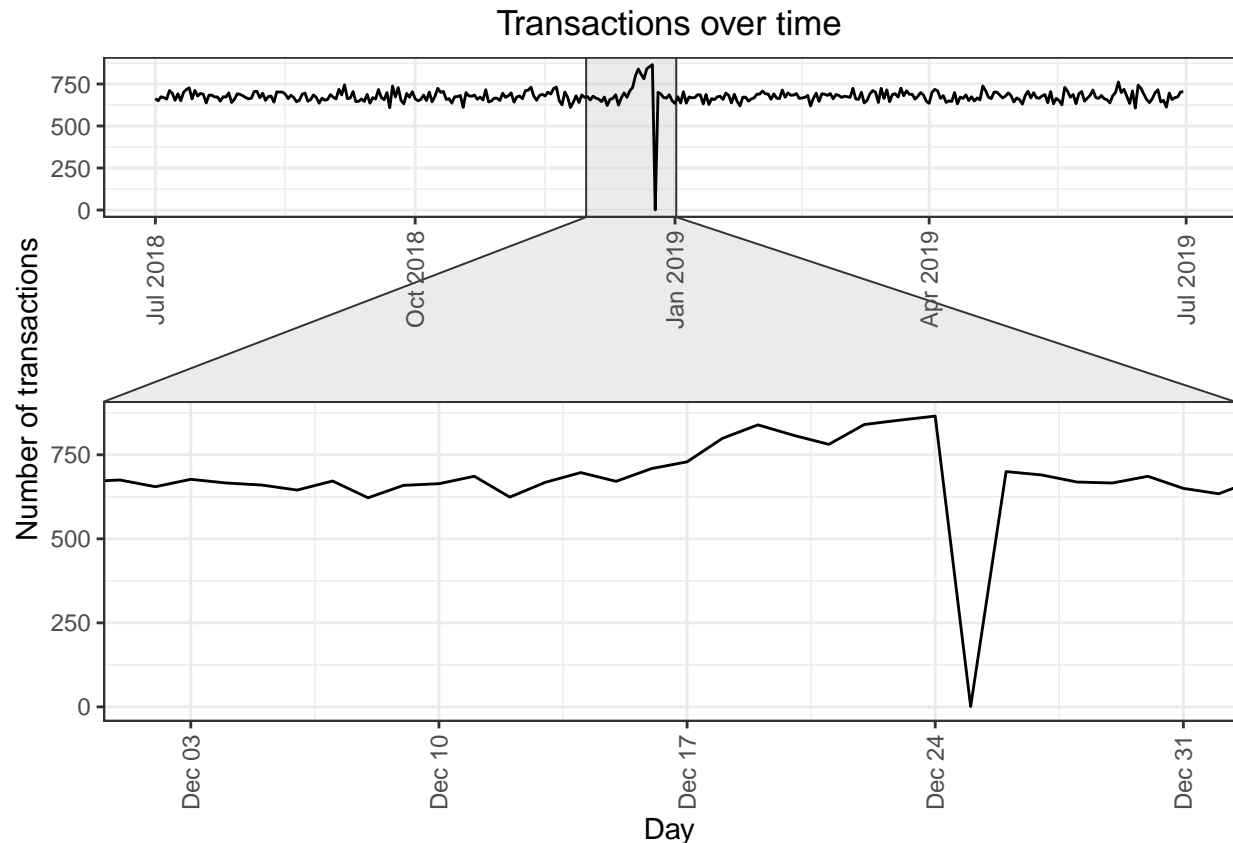
```
## Warning: `data_frame()` was deprecated in tibble 1.1.0.
## Please use `tibble()` instead.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was generated.
```

```
transactions_by_day <- merge(all_dates, transactionData, by = 'DATE', all.x=T) %>%
  group_by(DATE) %>%
  tally()
#### Setting plot themes to format graphs
theme_set(theme_bw())
theme_update(plot.title = element_text(hjust = 0.5))
#### Plot transactions over time
ggplot(transactions_by_day, aes(x = DATE, y = n)) +
  geom_line() +
  labs(x = "Day", y = "Number of transactions", title = "Transactions over time") +
  scale_x_date(breaks = "1 month") +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5))
```



We can see that there is an increase in purchases in December and a break in late December. Let's zoom in on this.

```
#### Filter to December and look at individual days
# Over to you - recreate the chart above zoomed in to the relevant dates.
ggplot(transactions_by_day, aes(x = DATE, y = n)) +
  geom_line() +
  facet_zoom(x= (DATE>as.Date('2018-12-01') & DATE<as.Date('2019-01-01')))) +
  labs(x = "Day", y = "Number of transactions", title = "Transactions over time") +
  # scale_x_date(breaks = "1 month") +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5))
```



We can see that the increase in sales occurs in the lead-up to Christmas and that there are zero sales on Christmas day itself. This is due to shops being closed on Christmas day.

Now that we are satisfied that the data no longer has outliers, we can move on to creating other features such as brand of chips or pack size from `PROD_NAME`. We will start with pack size.

```
#### Pack size
#### We can work this out by taking the digits that are in PROD_NAME
transactionData[, PACK_SIZE := parse_number(PROD_NAME)]

## Warning in `[.data.table`(transactionData, , `:=`(PACK_SIZE,
## parse_number(PROD_NAME)))': Invalid .internal.selfref detected and fixed by
## taking a (shallow) copy of the data.table so that := can add this new column
## by reference. At an earlier point, this data.table has been copied by R (or
## was created manually using structure() or similar). Avoid names<- and attr<-
## which in R currently (and oddly) may copy the whole data.table. Use set* syntax
## instead to avoid copying: ?set, ?setnames and ?setattr. If this message doesn't
## help, please report your use case to the data.table issue tracker so the root
## cause can be fixed or this message improved.

#### Always check your output
#### Let's check if the pack sizes look sensible
transactionData[, .N, PACK_SIZE][order(PACK_SIZE)]

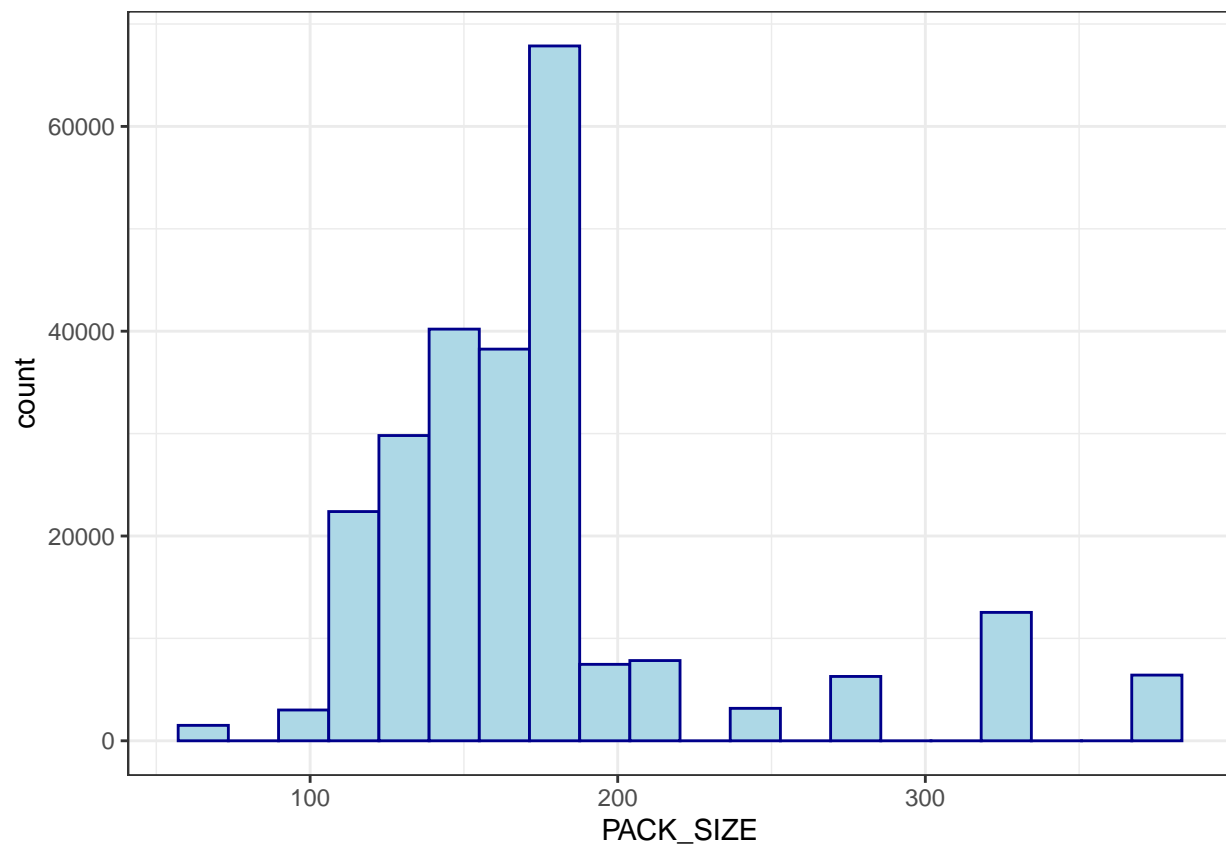
##   PACK_SIZE    N
## 1:      70 1507
## 2:      90 3008
## 3:     110 22387
## 4:     125 1454
```

```
## 5:      134 25102
## 6:      135 3257
## 7:      150 40203
## 8:      160 2970
## 9:      165 15297
## 10:     170 19983
## 11:     175 66390
## 12:     180 1468
## 13:     190 2995
## 14:     200 4473
## 15:     210 6272
## 16:     220 1564
## 17:     250 3169
## 18:     270 6285
## 19:     330 12540
## 20:     380 6416
```

The largest size is 380g and the smallest size is 70g - seems sensible!

*#### Let's plot a histogram of PACK\_SIZE since we know that it is a categorical variable and not a cont  
# Over to you! Plot a histogram showing the number of transactions by pack size.*

```
ggplot(transactionData, aes(PACK_SIZE)) +  
  geom_histogram(color = 'darkblue', fill = 'lightblue', bins = 20)
```



Pack sizes created look reasonable.

Now to create brands, we can use the first word in PROD\_NAME to work out the brand name...



```
#### Brands
# Over to you! Create a column which contains the brand of the product, by extracting it from the product name
#### Checking brands
# Over to you! Check the results look reasonable.
```

```
transactionData[, BRAND:= unlist(map(strsplit(transactionData$PROD_NAME, ' '),1))]
```

Some of the brand names look like they are of the same brands - such as RED and RRD, which are both Red Rock Deli chips. Let's combine these together.

```
#### Clean brand names
transactionData[BRAND == "Red", BRAND := "RRD"]

# Over to you! Add any additional brand adjustments you think may be required.
transactionData[BRAND == "Natural", BRAND := "NCC"]
transactionData[BRAND == "Smith", BRAND := "Smiths"]
transactionData[BRAND == "Grain", BRAND := "GrnWves"]
transactionData[BRAND == "Sunbites", BRAND := "Snbts"]
transactionData[BRAND == "Infuzions", BRAND := "Infzns"]
transactionData[BRAND == "Dorito", BRAND := "Doritos"]
#### Check again
# Over to you! Check the results look reasonable.
transactionData[, .N, , BRAND][order(desc(N))]
```

```
##      BRAND      N
## 1:   Kettle 41288
## 2:   Smiths 30353
## 3:   Doritos 25224
## 4: Pringles 25102
## 5:     RRD 16321
## 6:   Infzns 14201
## 7:    Thins 14075
## 8:      WW 10320
## 9:     Cobs 9693
##10: Tostitos 9471
##11: Twisties 9454
##12:   GrnWves 7740
##13:     NCC 7469
##14: Tyrrells 6442
##15: Cheezels 4603
##16:     CCs 4551
##17:   Snbts 3008
##18:   Cheetos 2927
##19:   Burger 1564
##20: Woolworths 1516
##21:   French 1418
##      BRAND      N
```

## Examining customer data

Now that we are happy with the transaction dataset, let's have a look at the customer dataset.

```
#### Examining customer data
# Over to you! Do some basic summaries of the dataset, including distributions of any key columns.
```

```
summary(customerData)
```

```
##  LYLTY_CARD_NBR      LIFESTAGE      PREMIUM_CUSTOMER  
##  Min.   :   1000    Length:72637      Length:72637  
##  1st Qu.:  66202    Class :character    Class :character  
##  Median : 134040    Mode  :character    Mode  :character  
##  Mean   : 136186  
##  3rd Qu.: 203375  
##  Max.   :2373711
```

```
customerData[, .N, PREMIUM_CUSTOMER][order(N)]
```

```
##  PREMIUM_CUSTOMER      N  
## 1:      Premium 18922  
## 2:      Budget 24470  
## 3:    Mainstream 29245
```

```
customerData[, .N, LIFESTAGE][order(N)]
```

```
##          LIFESTAGE      N  
## 1:    NEW FAMILIES  2549  
## 2: MIDAGE SINGLES/COUPLES 7275  
## 3:    YOUNG FAMILIES  9178  
## 4:    OLDER FAMILIES  9780  
## 5: YOUNG SINGLES/COUPLES 14441  
## 6: OLDER SINGLES/COUPLES 14609  
## 7:      RETIREES 14805
```

```
#### Merge transaction data to customer data
```

```
data <- merge(transactionData, customerData, all.x = TRUE)
```

As the number of rows in `data` is the same as that of `transactionData`, we can be sure that no duplicates were created. This is because we created `data` by setting `all.x = TRUE` (in other words, a left join) which means take all the rows in `transactionData` and find rows with matching values in shared columns and then joining the details in these rows to the `x` or the first mentioned table.

Let's also check if some customers were not matched on by checking for nulls.

```
# Over to you! See if any transactions did not have a matched customer.
```

```
sum(is.na(data))
```

```
## [1] 0
```

Great, there are no nulls! So all our customers in the transaction data has been accounted for in the customer dataset.

Note that if you are continuing with Task 2, you may want to retain this dataset which you can write out as a csv

```
fwrite(data, paste0(filePath, "QVI_data.csv"))
```

Data exploration is now complete!

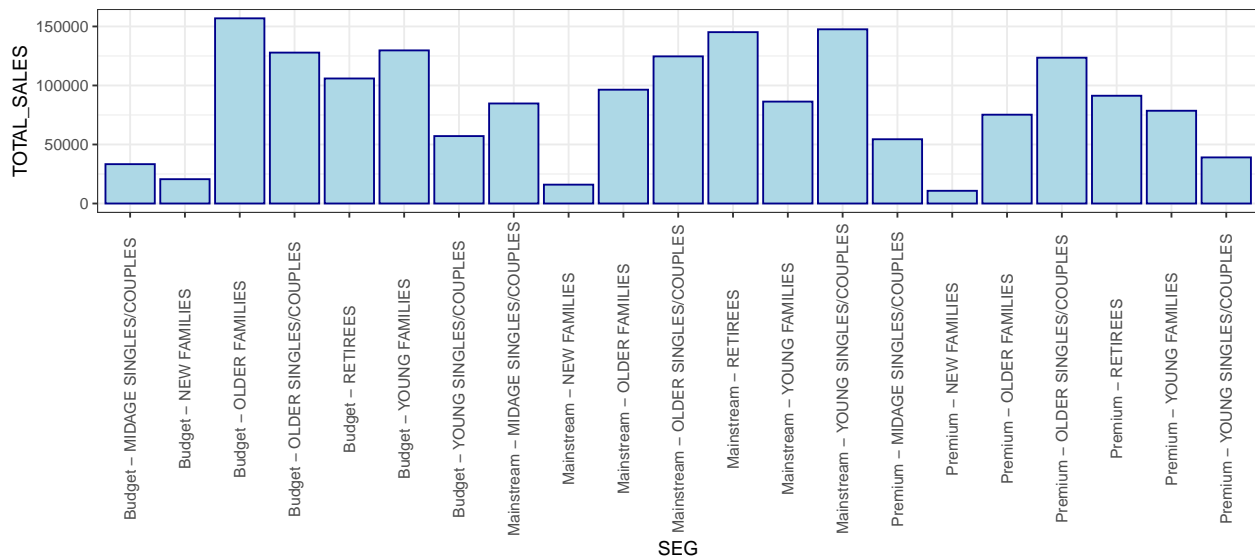
## Data analysis on customer segments

Now that the data is ready for analysis, we can define some metrics of interest to the client: - Who spends the most on chips (total sales), describing customers by lifestage and how premium their general purchasing behaviour is - How many customers are in each segment - How many chips are bought per customer by

segment - What's the average chip price by customer segment We could also ask our data team for more information. Examples are: - The customer's total spend over the period and total spend for each transaction to understand what proportion of their grocery spend is on chips - Proportion of customers in each customer segment overall to compare against the mix of customers who purchase chips

Let's start with calculating total sales by LIFESTAGE and PREMIUM\_CUSTOMER and plotting the split by these segments to describe which customer segment contribute most to chip sales.

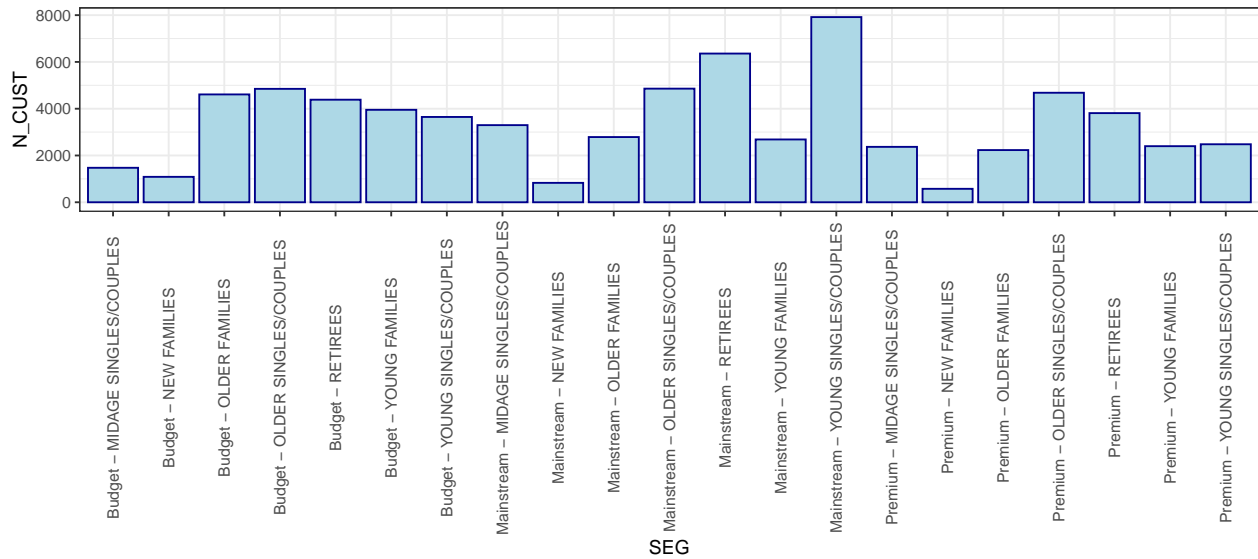
```
#### Total sales by LIFESTAGE and PREMIUM_CUSTOMER
# Over to you! Calculate the summary of sales by those dimensions and create a plot.
data %>%
  group_by(LIFESTAGE, PREMIUM_CUSTOMER) %>%
  summarise(TOTAL_SALES = sum(TOT_SALES), .groups = 'drop') %>%
  mutate(SEG = paste0(PREMIUM_CUSTOMER, ' - ', LIFESTAGE)) %>%
  ggplot(aes(SEG, TOTAL_SALES)) +
  geom_col(color = 'darkblue', fill = 'lightblue') +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5))
```



Sales are coming mainly from Budget - older families, Mainstream - young singles/couples, and Mainstream - retirees

Let's see if the higher sales are due to there being more customers who buy chips.

```
#### Number of customers by LIFESTAGE and PREMIUM_CUSTOMER
# Over to you! Calculate the summary of number of customers by those dimensions and create a plot.
data %>%
  group_by(LIFESTAGE, PREMIUM_CUSTOMER) %>%
  summarise(N_CUST = n_distinct(LYLT_CARD_NBR), .groups = 'drop') %>%
  mutate(SEG = paste0(PREMIUM_CUSTOMER, ' - ', LIFESTAGE)) %>%
  ggplot(aes(SEG, N_CUST)) +
  geom_col(color = 'darkblue', fill = 'lightblue') +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5))
```



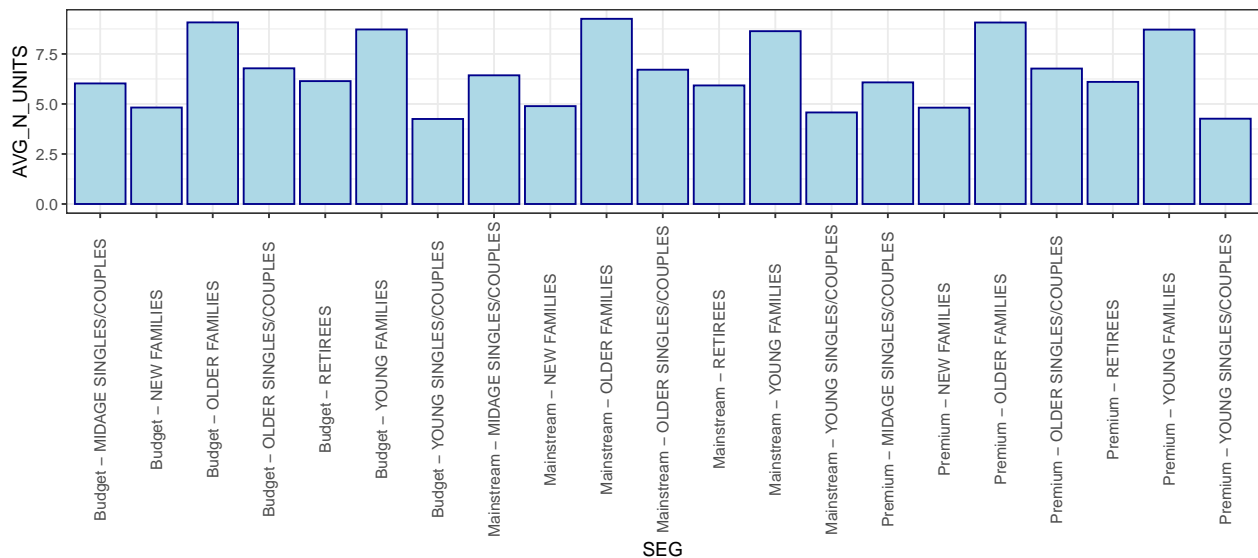
There are more Mainstream - young singles/couples and Mainstream - retirees who buy chips. This contributes to there being more sales to these customer segments but this is not a major driver for the Budget - Older families segment.

Higher sales may also be driven by more units of chips being bought per customer.

Let's have a look at this next.

#### Average number of units per customer by LIFESTAGE and PREMIUM\_CUSTOMER  
 # Over to you! Calculate and plot the average number of units per customer by those two dimensions.

```
data %>%
  group_by(LIFESTAGE, PREMIUM_CUSTOMER) %>%
  summarise(AVG_N_UNITS = sum(PROD_QTY)/n_distinct(LYLT_CARD_NBR),
            .groups = 'drop') %>%
  mutate(SEG = paste0(PREMIUM_CUSTOMER, ' - ', LIFESTAGE)) %>%
  ggplot(aes(SEG, AVG_N_UNITS, fill = factor(SEG)))+
  geom_col(color = 'darkblue', fill = 'lightblue')+
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5),
        legend.position = 'none')
```

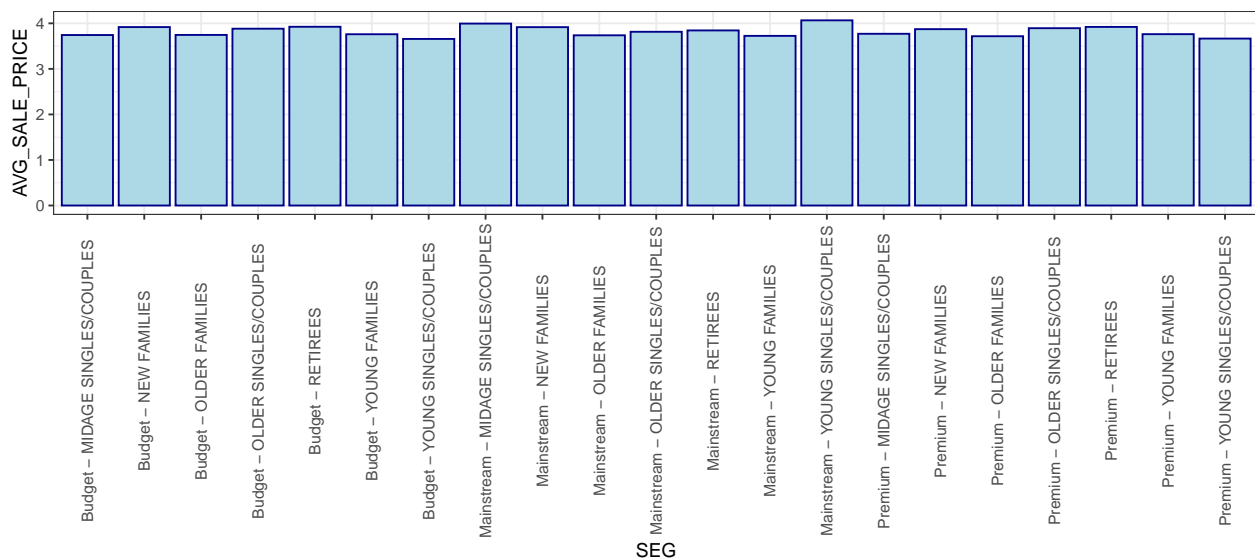


Older families and young families in general buy more chips per customer.

Let's also investigate the average price per unit chips bought for each customer segment as this is also a driver of total sales.

#### Average price per unit by LIFESTAGE and PREMIUM\_CUSTOMER  
 # Over to you! Calculate and plot the average price per unit sold (average sale price) by those two cus

```
avg_sale_price <- data %>%
  group_by(LIFESTAGE, PREMIUM_CUSTOMER) %>%
  summarise(AVG_SALE_PRICE = mean(TOT_SALES/PROD_QTY), .groups = 'drop') %>%
  mutate(SEG = paste0(PREMIUM_CUSTOMER, ' - ', LIFESTAGE))
avg_sale_price %>%
  ggplot(aes(SEG, AVG_SALE_PRICE, fill = SEG)) +
  geom_col(color = 'darkblue', fill = 'lightblue') +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5),
        legend.position = 'None')
```



Mainstream midage and young singles and couples are more willing to pay more per packet of chips compared to their budget and premium counterparts. This may be due to premium shoppers being more likely to buy healthy snacks and when they buy chips, this is mainly for entertainment purposes rather than their own consumption.

This is also supported by there being fewer premium midage and young singles and couples buying chips compared to their mainstream counterparts.

As the difference in average price per unit isn't large, we can check if this difference is statistically different.

#### Perform an independent t-test between mainstream vs premium and budget midage and  
 #### young singles and couples

# Over to you! Perform a t-test to see if the difference is significant.

```
YMSC <- data %>%
  mutate(UNIT_PRICE = TOT_SALES/PROD_QTY,
         MAINSTREAM = ifelse(PREMIUM_CUSTOMER == 'Mainstream', 'yes', 'no')) %>%
  filter(LIFESTAGE %in% c('YOUNG SINGLES/COUPLES', 'MIDAGE SINGLES/COUPLES'))
attach(YMSC)
# library(car)
# leveneTest(UNIT_PRICE ~ MAINSTREAM)
test = t.test(UNIT_PRICE[MAINSTREAM=='yes'], UNIT_PRICE[MAINSTREAM=='no'],
```

```

alternative = 'greater')
test

##
## Welch Two Sample t-test
##
## data: UNIT_PRICE[MAINSTREAM == "yes"] and UNIT_PRICE[MAINSTREAM == "no"]
## t = 37.624, df = 54791, p-value < 2.2e-16
## alternative hypothesis: true difference in means is greater than 0
## 95 percent confidence interval:
## 0.3187234 Inf
## sample estimates:
## mean of x mean of y
## 4.039786 3.706491

test$p.value

## [1] 3.483677e-306

detach(YMSC)

```

The t-test results in a p-value of 3.483677e-306, i.e. the unit price for mainstream, young and mid-age singles and couples [ARE] significantly higher than that of budget or premium, young and midage singles and couples.

## Deep dive into specific customer segments for insights

We have found quite a few interesting insights that we can dive deeper into.

We might want to target customer segments that contribute the most to sales to retain them or further increase sales. Let's look at Mainstream - young singles/couples. For instance, let's find out if they tend to buy a particular brand of chips.

```

#### Deep dive into Mainstream, young singles/couples
# Over to you! Work out if there are brands that these two customer segments prefer more than others.
# You could use a technique called affinity analysis or a-priori analysis (or any other method if you p
target_SEG <- data %>%
  filter(
    PREMIUM_CUSTOMER == 'Mainstream' &
    LIFESTAGE == 'YOUNG SINGLES/COUPLES'
  )
other_SEG <- data %>%
  filter(!TXN_ID %in% target_SEG$TXN_ID)
target_transaction <- plyr::ddply(target_SEG, c('LYLTY_CARD_NBR'),
  function(df1)paste(unique(df1$BRAND), collapse = ','))
target_transaction <- target_transaction %>%
  select(-LYLTY_CARD_NBR) %>%
  transmute(items = factor(V1))
fwrite(target_transaction, paste0(filePath,"tr_target_brand.csv"))

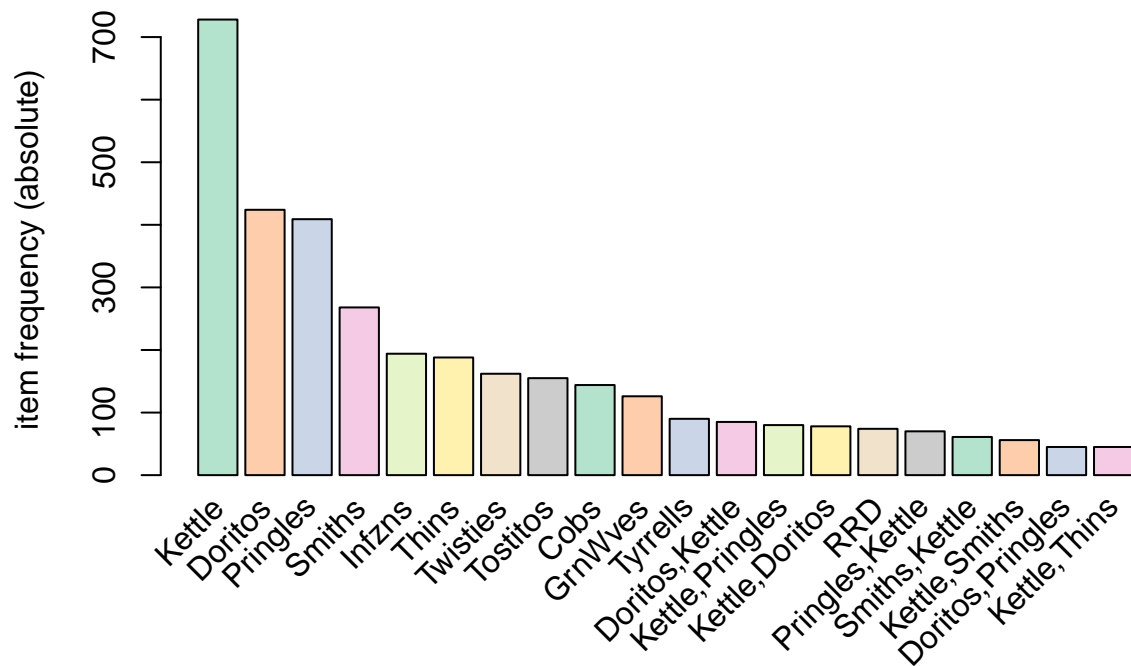
other_transaction <- plyr::ddply(other_SEG, c("LYLTY_CARD_NBR"),
  function(df1)paste(unique(df1$BRAND), collapse = ','))
other_transaction <- other_transaction %>%
  select(-LYLTY_CARD_NBR) %>%
  transmute(items = factor(V1))
fwrite(other_transaction, paste0(filePath,"tr_other_brand.csv"))

```

```
tr_target <- read.transactions(paste0(filePath,"tr_target_brand.csv"),
                              format = 'basket', sep = ',')
tr_other <- read.transactions(paste0(filePath, "tr_other_brand.csv"),
                              format = 'basket', sep = ',')

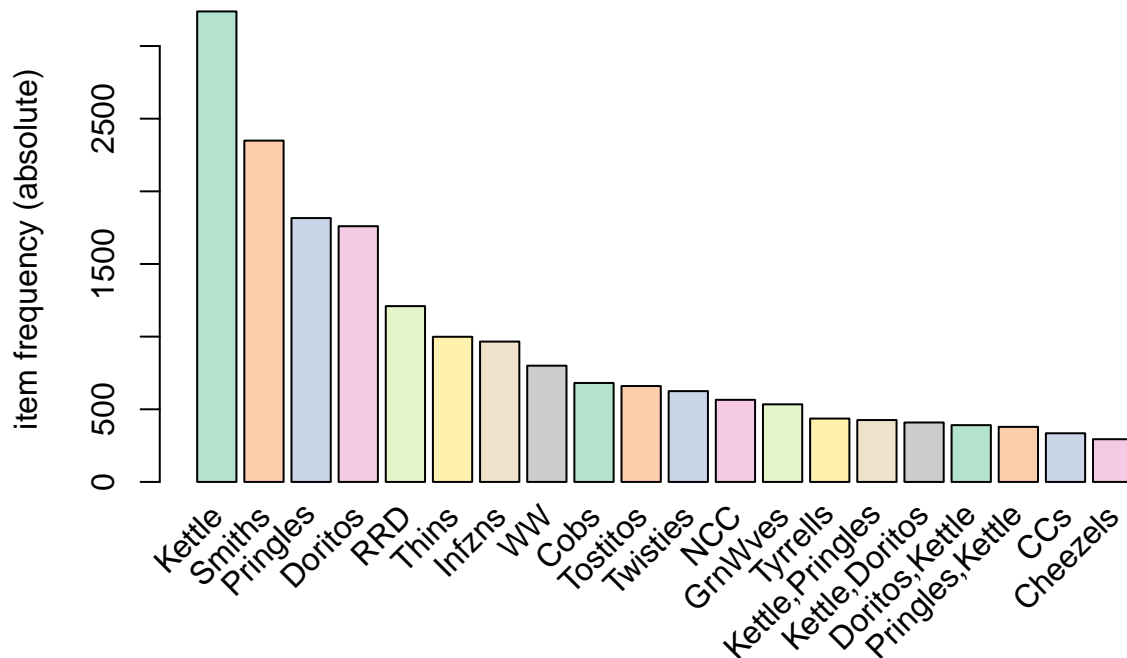
p1 <- itemFrequencyPlot(tr_target, topN = 20, type = 'absolute',
                        col = brewer.pal(8, 'Pastel2'),
                        main= 'Item Frequency - Target')
```

### Item Frequency – Target



```
p2 <- itemFrequencyPlot(tr_other, topN = 20, type = 'absolute',
                        col = brewer.pal(8, 'Pastel2'),
                        main = 'Item Frequency - Others')
```

## Item Frequency – Others



We can see that:

1. Kettle is the most popular brand for both our target customer segment(Mainstream - YOUNG SINGLES/COUPLES) and other segments.
2. Doritos, Pringles, and Smiths seems popular for all segments while Smiths is not as popular for our target segment and Doritos not as popular for other segments.
3. Twisties seems to gain some attention within our target segments ranking at No.7 among all itemsets.
4. RRD is not as favorable to our target segments compared to other segments.

Let's also find out if our target segment tends to buy larger packs of chips.

*#### Preferred pack size compared to the rest of the population  
# Over to you! Do the same for pack size.*

```
target_transaction <- plyr::ddply(target_SEG, c('LYLTY_CARD_NBR', "DATE"),
                                function(df1)paste(unique(df1$PACK_SIZE), collapse = ','))
target_transaction <- target_transaction %>%
  select(-c(LYLYTY_CARD_NBR,DATE)) %>%
  transmute(items = factor(V1))
fwrite(target_transaction, paste0(filePath,"tr_target_pack_size.csv"))

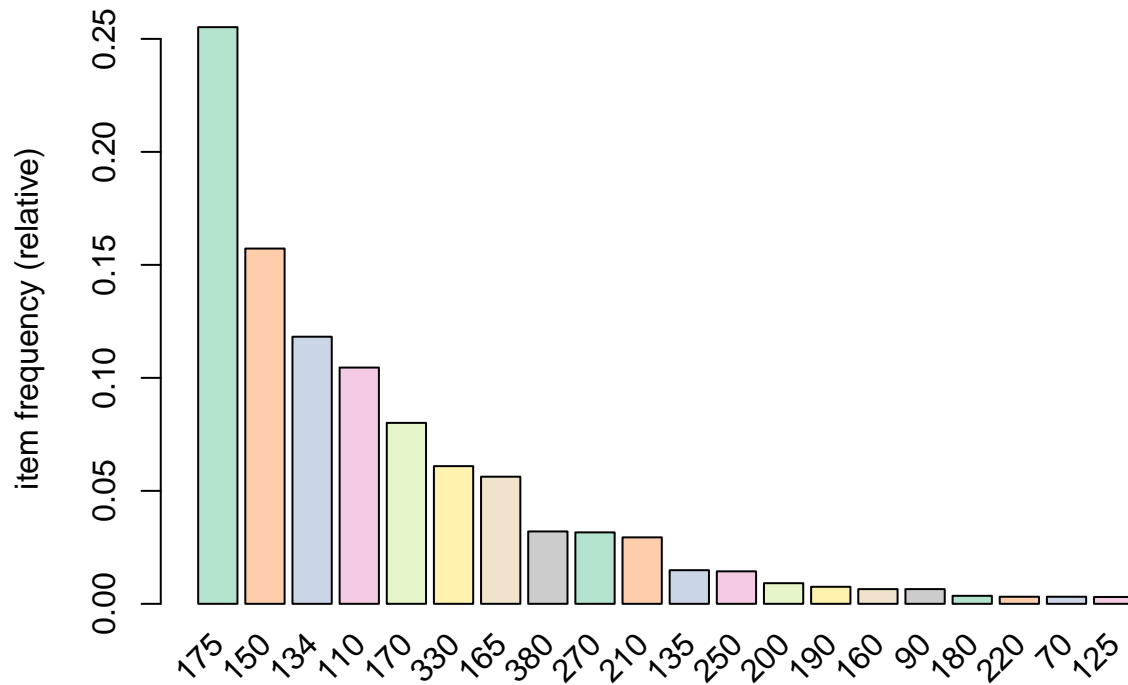
other_transaction <- plyr::ddply(other_SEG, c("LYLTY_CARD_NBR", "DATE"),
                                function(df1)paste(unique(df1$PACK_SIZE), collapse = ','))
other_transaction <- other_transaction %>%
  select(-c(LYLYTY_CARD_NBR, DATE)) %>%
  transmute(items = factor(V1))
fwrite(other_transaction, paste0(filePath,"tr_other_pack_size.csv"))

tr_target <- read.transactions(paste0(filePath,"tr_target_pack_size.csv"),
                              format = 'basket', sep = ',')
tr_other <- read.transactions(paste0(filePath, "tr_other_pack_size.csv"),
                              format = 'basket', sep = ',')
```



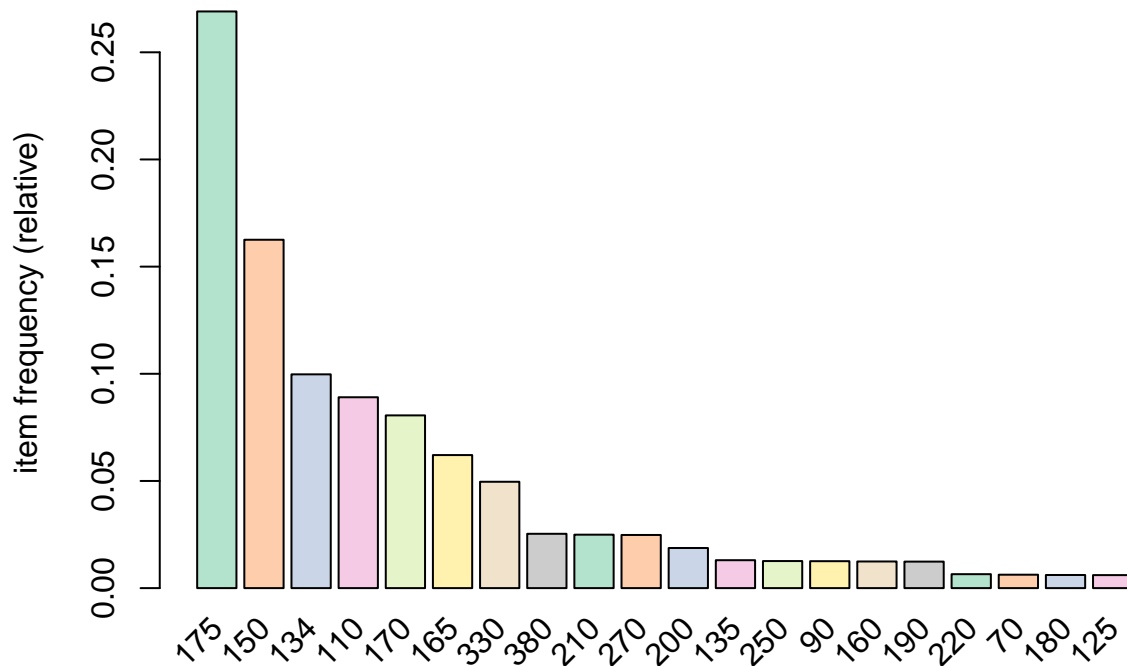
```
p1 <- itemFrequencyPlot(tr_target, topN = 20, type = 'relative',
  col = brewer.pal(8, 'Pastel2'),
  main= 'Item Frequency - Target')
```

## Item Frequency – Target



```
p2 <- itemFrequencyPlot(tr_other, topN = 20, type = 'relative',
  col = brewer.pal(8, 'Pastel2'),
  main = 'Item Frequency - Others')
```

## Item Frequency – Others



Generally, the distribution are similar in the item frequency plot, indicating that tendency is unclear in this plot for target segment buying larger packs of chips compared to other segments. However, we can still see our target segment in favor of 330g over 165g for 6th and 7th most popular pack sizes, which differ from the other segments. Cases are the same for the 9th and 10th most popular pack sizes (270g and 210g). Thus, we conduct a t-test for pack size over target segments and other segments:

```
target_data = data %>%
  mutate(TARGET = ifelse(PREMIUM_CUSTOMER == 'Mainstream' &
    LIFESTAGE == 'YOUNG SINGLES/COUPLES',
    'yes', 'no'))
test = t.test(PACK_SIZE ~ TARGET, data = target_data,
  alternative = 'less')
test

##
## Welch Two Sample t-test
##
## data: PACK_SIZE by TARGET
## t = -6.3297, df = 22505, p-value = 1.251e-10
## alternative hypothesis: true difference in means between group no and group
yes is less than 0
## 95 percent confidence interval:
## -Inf -2.219061
## sample estimates:
## mean in group no mean in group yes
## 175.3460 178.3442

paste0('p-value: ', test$p.value)

## [1] "p-value: 1.25110481330845e-10"
```

The t-test results in a p-value of 1.251105e-10. i.e. the pack size of Mainstream - YOUNG SINGLES/COUPLES

are significantly higher than the other segments.