# Task1- Data preparation and customer analytics

SIDDHESH POTE

14/12/2020

**Loading the required libraries**

```
pacman::p_load(ggplot2, dplyr , readxl ,data.table, ggmosaic, readr)
```

**storing the data in the required variables**

```
# Point the filePath to where you have downloaded the datasets to and assign the data files to data.tab
filepath <- "C:/Users/Siddhesha/Desktop/R commom directory/quantinum virtual internship/"

transactionData <- read_xlsx("C:/Users/Siddhesha/Desktop/R commom directory/quantinum virtual internship
customerData <-  fread(paste0(filepath,"QVI_purchase_behaviour.csv"))
transactionData <- data.table(transactionData)
```

## Exploratory data analyis

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

```
str(transactionData)
```

```
## Classes 'data.table' and 'data.frame':   264836 obs. of  8 variables:
##  $ DATE          : num  43390 43599 43605 43329 43330 ...
##  $ STORE_NBR     : num  1 1 1 2 2 4 4 4 5 7 ...
##  $ LYLTY_CARD_NBR: num  1000 1307 1343 2373 2426 ...
##  $ TXN_ID        : num  1 348 383 974 1038 ...
##  $ PROD_NBR      : num  5 66 61 69 108 57 16 24 42 52 ...
##  $ PROD_NAME     : chr  "Natural Chip        Compny SeaSalt175g" "CCs Nacho Cheese    175g" "Smiths (
##  $ PROD_QTY      : num  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>
```

Here we can see that the all the variables are in their correct format i.e the integer format and product name is in character format except for the date variable which should have been in the date format is in the form of integer format... The date starts from 43282 i.e the no of days... in excel the dates are recorded form 30 DEC 1899... That means our data starts at the date that is 43282 days away from the origin date.. hence we convert the date variable in the form of date

### examining the DATE variable

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

now we can see that the dates are in their respective form

## examining the PROD_NAME variable

```
transactionData[, .N, PROD_NAME]
```

```
##                                  PROD_NAME    N
##   1:    Natural Chip        Compny SeaSalt175g 1468
##   2:                  CCs Nacho Cheese    175g 1498
##   3:    Smiths Crinkle Cut  Chips Chicken 170g 1484
##   4:    Smiths Chip Thinly  S/Cream&Onion 175g 1473
##   5: Kettle Tortilla ChpsHny&Jlpno Chili 150g 3296
##  ---
## 110:    Red Rock Deli Chikn&Garlic Aioli 150g 1434
## 111:      RRD SR Slow Rst     Pork Belly 150g 1526
## 112:                 RRD Pc Sea Salt     165g 1431
## 113:       Smith Crinkle Cut   Bolognese 150g 1451
## 114:                 Doritos Salsa Mild  300g 1472
```

the data shows that there are 114 types of chips/salsa/rings etc product that were sold... since we are only interested in the potato chips data we would like to keep only the data of potato ships and discard other. We can do some basic text analysis by summarising the individual words in the product name.

```
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().

```
# Removing digits
productWords <- productWords[grepl("\\d", words) == FALSE, ]

# Removing special characters
productWords <- productWords[grepl("[:alpha:]", words), ]

# Let's look at the most common words by counting the number of times a word appears and sorting them by
productWords[, .N, words][order(N, decreasing = TRUE)]
```

```
##          words  N
##   1:     Chips 21
##   2:    Smiths 16
##   3:   Crinkle 14
##   4:    Kettle 13
##   5:    Cheese 12
```

```
##   ---
## 127: Chikn&Garlic   1
## 128:        Aioli   1
## 129:         Slow   1
## 130:        Belly   1
## 131:    Bolognese   1
```

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

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

# summarizing the data
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
```

In product quantity it is visible that in PROD_QTY column the max value is of 200 i.e 200 quantity were purchased at once.. hence this can be considered as a outlier to our data.. there are no nulls as all the summary statistics have a numerical value # finding the outlier in PROD_QTY

```
# Filter the dataset to find the outlier
transactionData[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
```

```
# as we can that there were two transactions done that had product quantity > 10. Both the transaction i

# Let's see if the customer has had other transactions
transactionData[LYLTY_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
```

```
# No other transactions were done by the cust except for those two where he purchased 200 qty of chips.

# Filter out the customer based on the loyalty card number
transactionData <- transactionData[LYLTY_CARD_NBR != 226000, ]
#### Re-examine transaction data
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.

```
## exploring the number of transactions according to the date

transactionData[, .N, by = DATE]
```
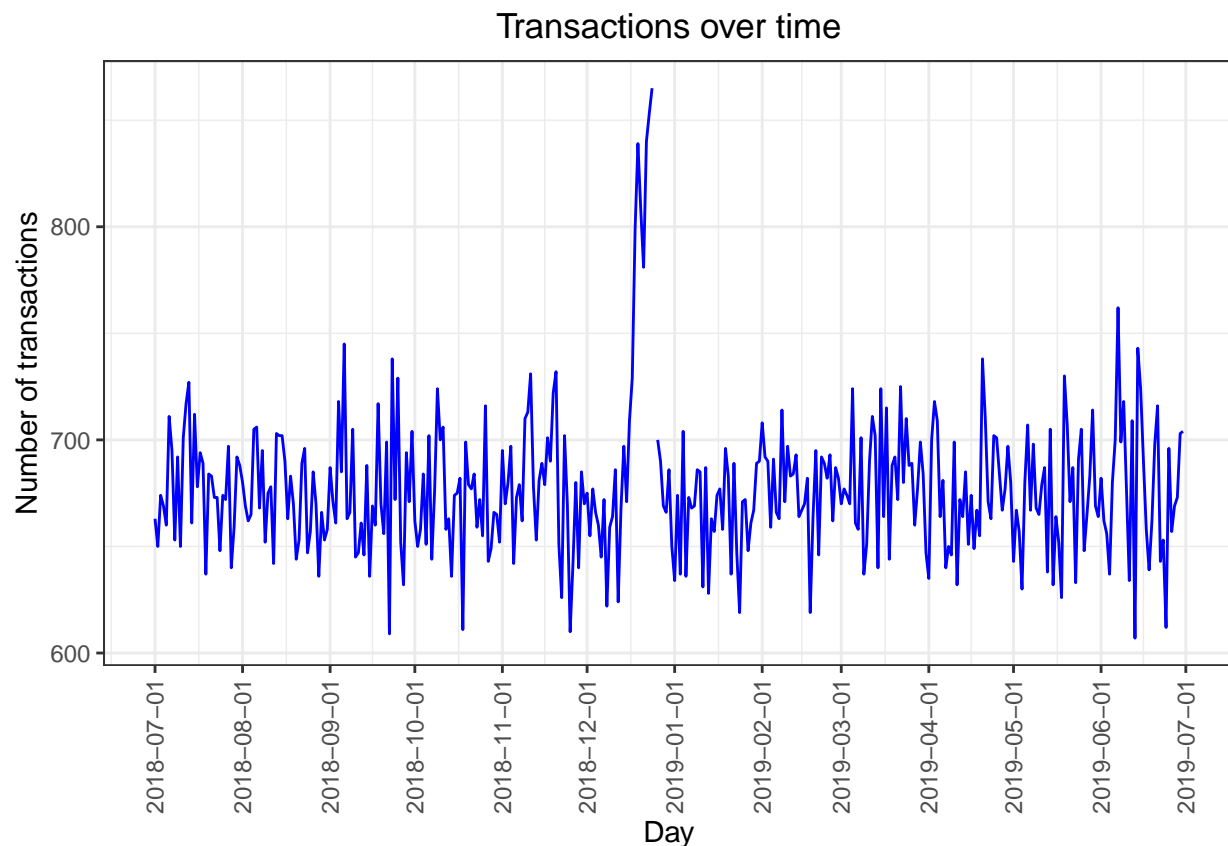
```
##            DATE   N
##   1: 2018-10-17 682
##   2: 2019-05-14 705
##   3: 2019-05-20 707
##   4: 2018-08-17 663
##   5: 2018-08-18 683
##  ---
## 360: 2018-12-08 622
## 361: 2019-01-30 689
## 362: 2019-02-09 671
## 363: 2018-08-31 658
## 364: 2019-02-12 684
```

Here as we can see that there are only 364 rows and not 365. . . the frequency of transaction as per date is shown in frequency column i.e 663 trans were done on 1 july 2018 and so on. . . . as the data is for a year so there should have been 365 rows and not 364. . . hence we check for the missing data
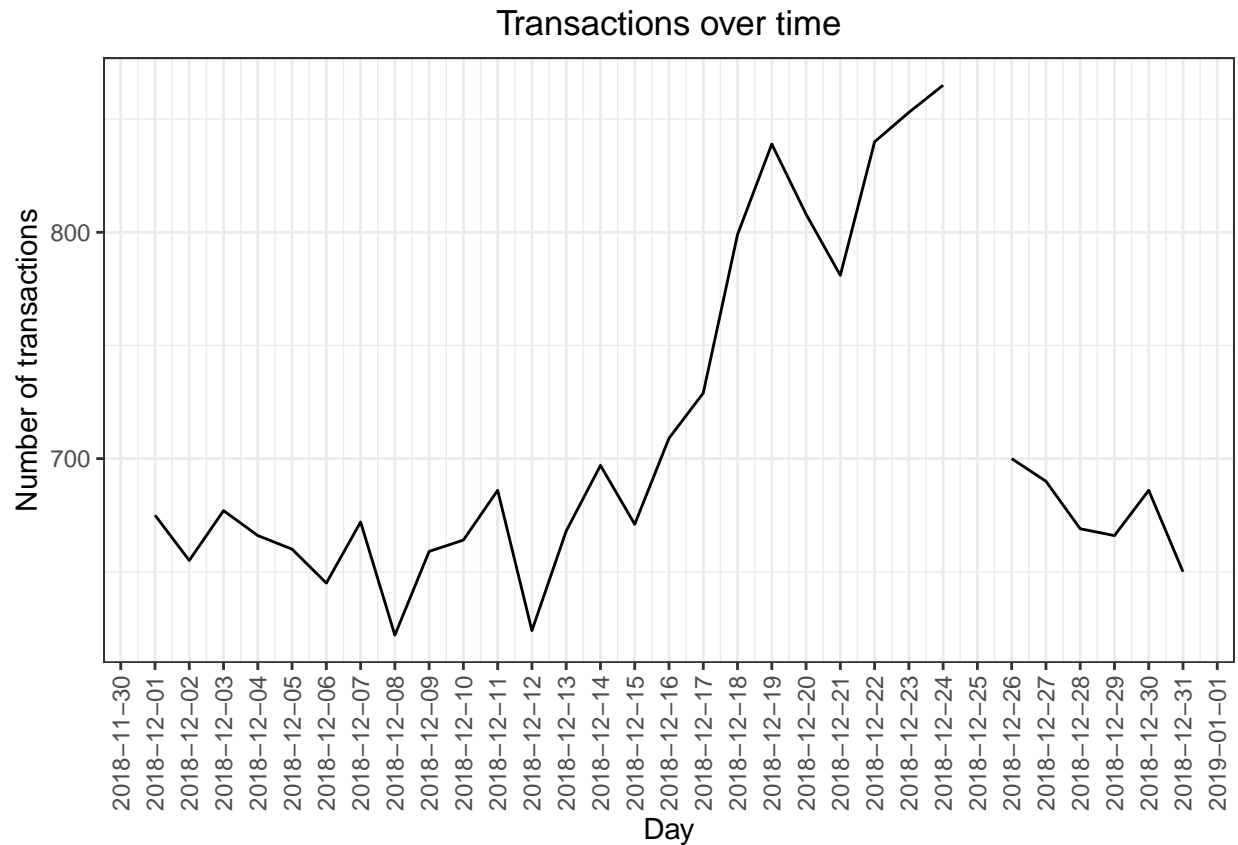
**creating a plot to check the missing date**

```r
# Create a sequence of dates and join this the count of transactions by date
allDates <- data.table(seq(as.Date("2018/07/01"), as.Date("2019/06/30"), by = "day"))
setnames(allDates, "DATE")
transactions_by_day <- merge(allDates, transactionData[, .N, by = DATE], all.x = TRUE)
### 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(col = "blue") +
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.

```r
# Filter to December and look at individual days
ggplot(transactions_by_day[month(DATE) == 12, ], aes(x = DATE, y = N)) +
geom_line() +
labs(x = "Day", y = "Number of transactions", title = "Transactions over time") +
scale_x_date(breaks = "1 day") +
theme(axis.text.x = element_text(angle = 90, vjust = 0.5))
```

## Transactions over time



From this we can conclude that there was an increase in sale due to christmas i.e 25th december and because of shop being closed on that day no data was recorded and the data was recorded for the day... hence we are satisfied that we don't have any other missing value and we can proceed further to create 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)]

# 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!

we created a new column that contains the packsize of the product.. lets see that the pack sizes are not too big or small... i.e they are within a specific range or not
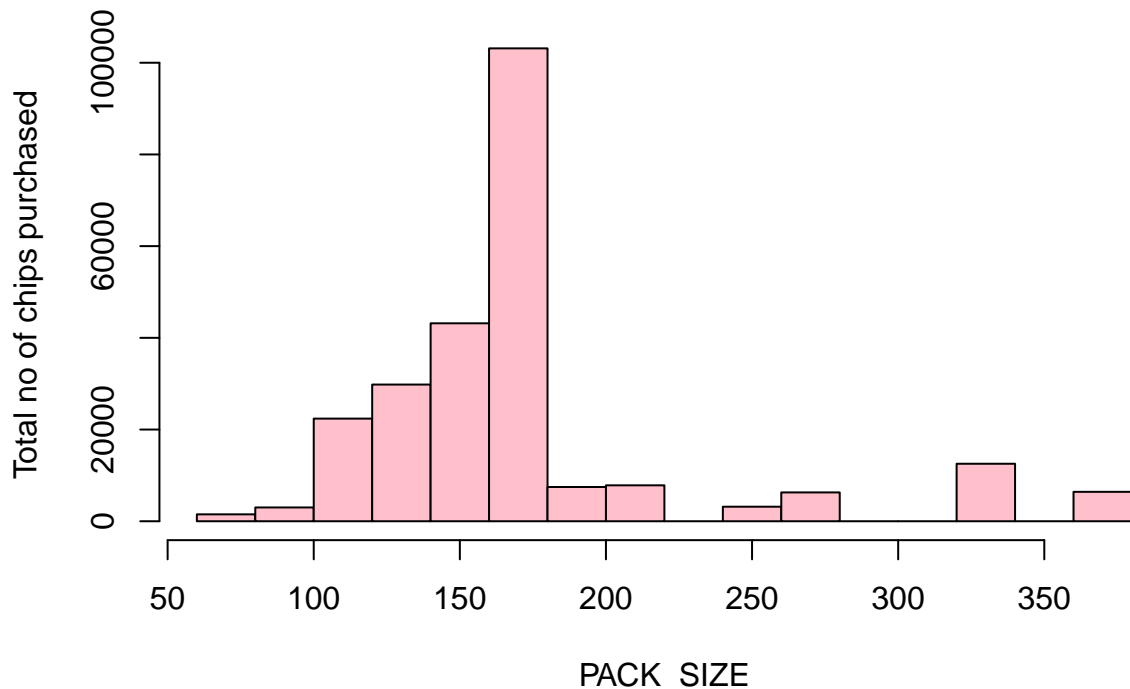
```
# Let's check the output of the first few rows to see if we have indeed picked out pack size.
transactionData
```

```
##               DATE STORE_NBR LYLTY_CARD_NBR TXN_ID PROD_NBR
##       1: 2018-10-17         1           1000      1        5
##       2: 2019-05-14         1           1307    348       66
##       3: 2019-05-20         1           1343    383       61
##       4: 2018-08-17         2           2373    974       69
##       5: 2018-08-18         2           2426   1038      108
##      ---
## 246736: 2019-03-09       272         272319 270088       89
## 246737: 2018-08-13       272         272358 270154       74
## 246738: 2018-11-06       272         272379 270187       51
## 246739: 2018-12-27       272         272379 270188       42
## 246740: 2018-09-22       272         272380 270189       74
##                                        PROD_NAME PROD_QTY TOT_SALES PACK_SIZE
##       1:    Natural Chip        Compny SeaSalt175g        2       6.0       175
##       2:                  CCs Nacho Cheese    175g        3       6.3       175
##       3:    Smiths Crinkle Cut  Chips Chicken 170g        2       2.9       170
##       4:    Smiths Chip Thinly  S/Cream&Onion 175g        5      15.0       175
##       5: Kettle Tortilla ChpsHny&Jlpno Chili 150g        3      13.8       150
##      ---
## 246736:  Kettle Sweet Chilli And Sour Cream 175g        2      10.8       175
## 246737:             Tostitos Splash Of  Lime 175g        1       4.4       175
## 246738:                 Doritos Mexicana    170g        2       8.8       170
## 246739:  Doritos Corn Chip Mexican Jalapeno 150g        2       7.8       150
## 246740:             Tostitos Splash Of  Lime 175g        2       8.8       175
```

Lets plot a histogram of PACK_SIZE since we know its a categorial variable

```
## plotting histogram of the packsize
options(scipen=999) # turn off scientific notations like 1e+05
hist(transactionData[, PACK_SIZE], col = "pink",border = "black" , xlab = "PACK  SIZE", ylab = "Total n
```

# TOGRAM OF NO. OF CHIPS PURCHASED ACCORDING TO THEIR PACK



PACK SIZE

the plot looks reasonable with no outliers and from the plot it can be seen that the the packs of size 170-180
was purchased the most

## Brands

```
# we'll use the first word of the PROD_NAME to create our data of brand

transactionData[, BRAND := toupper(substr(PROD_NAME, 1, regexpr(pattern = ' ', PROD_NAME) - 1))]

# Checking brands
transactionData[, .N, by = BRAND][order(-N)]
```

```
##            BRAND     N
##  1:       KETTLE 41288
##  2:       SMITHS 27390
##  3:     PRINGLES 25102
##  4:      DORITOS 22041
##  5:        THINS 14075
##  6:          RRD 11894
##  7:     INFUZIONS 11057
##  8:           WW 10320
##  9:         COBS  9693
## 10:      TOSTITOS  9471
## 11:      TWISTIES  9454
```

```
## 12:    TYRRELLS  6442
## 13:       GRAIN  6272
## 14:     NATURAL  6050
## 15:    CHEEZELS  4603
## 16:         CCS  4551
## 17:         RED  4427
## 18:      DORITO  3183
## 19:       INFZNS  3144
## 20:       SMITH  2963
## 21:     CHEETOS  2927
## 22:       SNBTS  1576
## 23:      BURGER  1564
## 24: WOOLWORTHS  1516
## 25:     GRNWVES  1468
## 26:    SUNBITES  1432
## 27:         NCC  1419
## 28:      FRENCH  1418
##           BRAND     N
```

Here we can see that chips of kettle brand have been purchased the most... also the data has no outliers in it.. the only problem is with the brand name red and RRD which both are same... hence we need to merge the data which contains rrd and red as the brand name

```
transactionData[BRAND == "RED", BRAND := "RRD"]
transactionData[BRAND == "SNBTS", BRAND := "SUNBITES"]
transactionData[BRAND == "INFZNS", BRAND := "INFUZIONS"]
transactionData[BRAND == "WW", BRAND := "WOOLWORTHS"]
transactionData[BRAND == "SMITH", BRAND := "SMITHS"]
transactionData[BRAND == "NCC", BRAND := "NATURAL"]
transactionData[BRAND == "DORITO", BRAND := "DORITOS"]
transactionData[BRAND == "GRAIN", BRAND := "GRNWVES"]

# Check again
transactionData[, .N, by = BRAND][order(BRAND)]
```

```
##            BRAND     N
##  1:       BURGER  1564
##  2:          CCS  4551
##  3:      CHEETOS  2927
##  4:     CHEEZELS  4603
##  5:         COBS  9693
##  6:      DORITOS 25224
##  7:       FRENCH  1418
##  8:      GRNWVES  7740
##  9:    INFUZIONS 14201
## 10:       KETTLE 41288
## 11:      NATURAL  7469
```

```
## 12:    PRINGLES 25102
## 13:         RRD 16321
## 14:      SMITHS 30353
## 15:    SUNBITES  3008
## 16:       THINS 14075
## 17:    TOSTITOS  9471
## 18:    TWISTIES  9454
## 19:    TYRRELLS  6442
## 20: WOOLWORTHS 11836
```

```
# now 8 of our rows that had similar brand has been merged and now we can finally stop our data explora
```

## Checking the customer data

```
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 SI
##  $ PREMIUM_CUSTOMER: chr  "Premium" "Mainstream" "Budget" "Mainstream" ...
##  - attr(*, ".internal.selfref")=<externalptr>
```

```
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
```

we can see that the data is mostly descriptive. i,e it gives the description of the customer who purchased the chips... the loyalty card number is a numeric vector while lifestage and premium_customer are character vectors

Let's have a closer look at the LIFESTAGE and PREMIUM_CUSTOMER columns.

```
## Examining the values of lifestage and premium_customer
customerData[, .N, by = LIFESTAGE][order(-N)]
```

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

```
customerData[, .N, by = PREMIUM_CUSTOMER][order(-N)]
```

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

As there do not seem to be any issues with the customer data, we can now go ahead and join the transaction and customer data sets together

```
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.

```
colSums(is.na(data))
```

```
##    LYLTY_CARD_NBR            DATE       STORE_NBR             TXN_ID
##                 0               0               0                  0
##          PROD_NBR       PROD_NAME        PROD_QTY          TOT_SALES
##                 0               0               0                  0
##         PACK_SIZE           BRAND       LIFESTAGE PREMIUM_CUSTOMER
##                 0               0               0                  0
```

here we can see in all columns there are no NA values hence we can proceed further since all the data was matched properly

```
# saving the cleaned file for further analysis
# write.csv(data, "Cleaned_data.csv")
```

Our data exploration is now complete...

## Performing data analysis on customer segments

since the data is ready for data analysis we can now create various questions and define our interest on the variable of interest such as. - 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
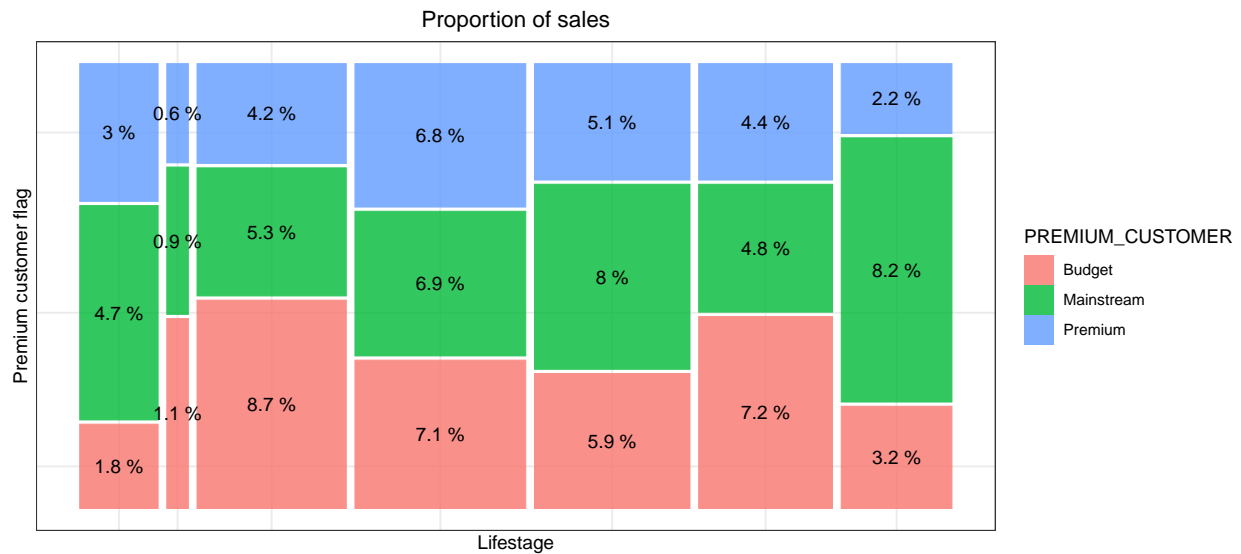
Let's start with calculating total sales by LIFESTAGE and PREMIUM_CUSTOMER and plotting the split by these segments to describe which customer segment contribute most to chip sales.

```
# Total sales by LIFESTAGE and PREMIUM_CUSTOMER
sales <- data[, .(SALES = sum(TOT_SALES)), .(LIFESTAGE,PREMIUM_CUSTOMER)]
# create plot
p <- ggplot(data = sales) +
  geom_mosaic(aes(weight = SALES, x = product(PREMIUM_CUSTOMER, LIFESTAGE) , fill = PREMIUM_CUSTOMER))
```

```
  labs(x = "Lifestage", y = "Premium customer flag", title = "Proportion of sales") + theme(axis.text.x

# Plot and label with proportion of sales
p +
  geom_text(data = ggplot_build(p)$data[[1]], aes(x = (xmin + xmax)/2 , y = (ymin + ymax)/2, label = as
```
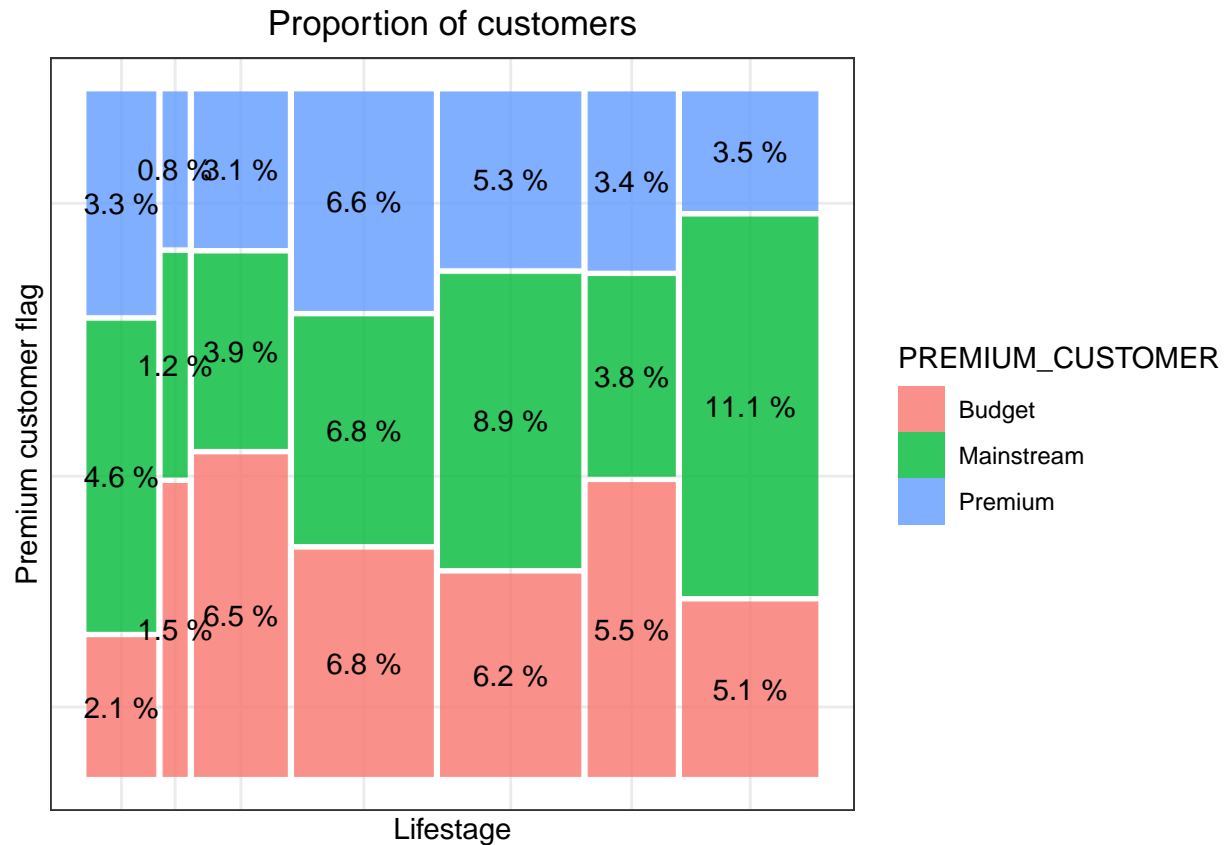


we can see from the plot that the sales are mostly due to the budget- older families, mainstream young single/couples and mainstream - retirees

lets see if the highers sales are due to there being more customers who buy chips..

```
## Number of customers by LIFESTAGE and PREMIUM_CUSTOMER
customers <- data[, .(CUSTOMERS = uniqueN(LYLTY_CARD_NBR)), .(LIFESTAGE, PREMIUM_CUSTOMER)][order(-CUST
labels <- c("A", "b", "c", "D", "e", "f")
# Create plot
p <- ggplot(data = customers) + geom_mosaic(aes(weight = CUSTOMERS, x = product(PREMIUM_CUSTOMER, LIFES

p + geom_text(data = ggplot_build(p)$data[[1]], aes(x = (xmin + xmax)/2 , y = (ymin + ymax)/2, label =
```
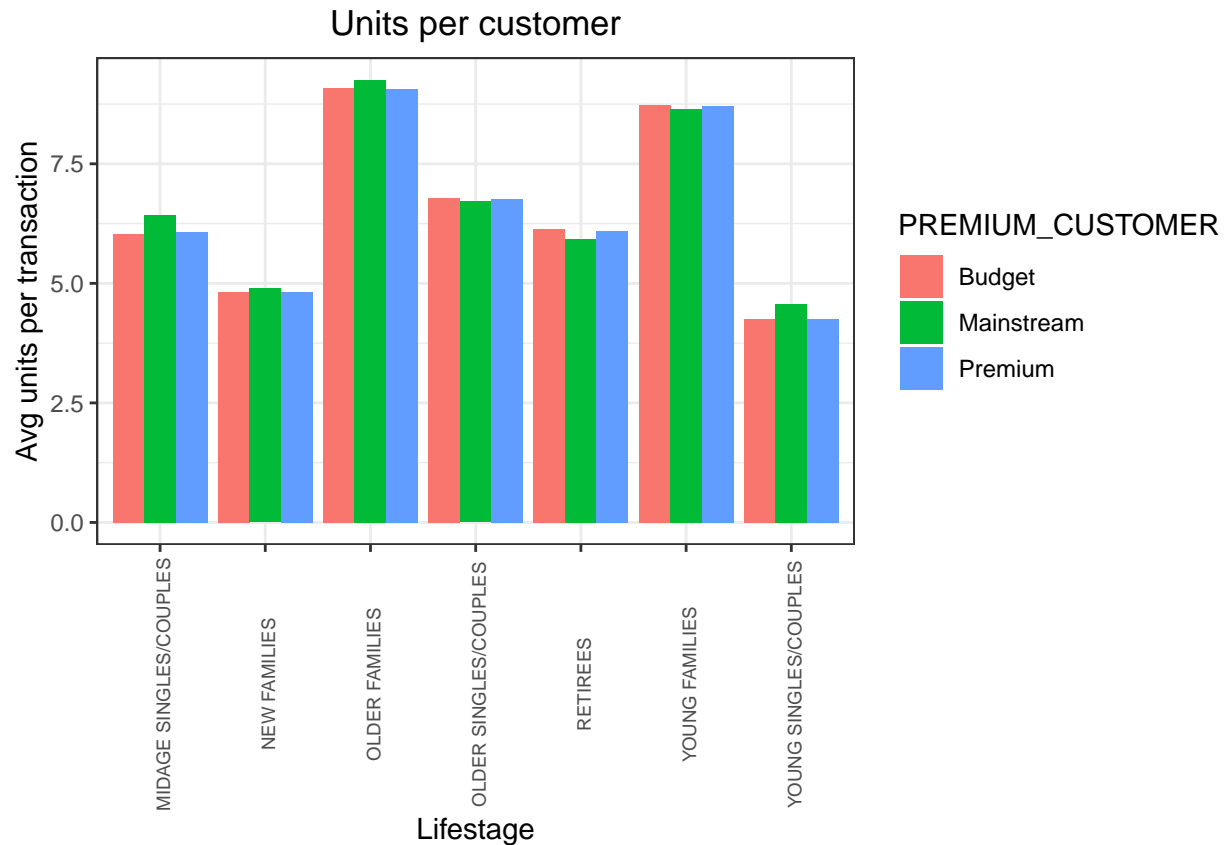
## Proportion of customers



From here we can see that mainstream- young/single couples and mainstream retirees contribute most to the sales of chips but it is not a major driver for budget- older families segment.

Higher sales may also be driven by no of chips bought by each customer.. hence we'll try to plot average no of chips i.e average no of PROD_QTY by lifestage and premium_customer

Higher sales may also be driven by more units of chips being bought per customer. Let's have a look at this next

```r
# Finding the average quantity of chips bought by each customers
avg_units <- data[, .(AVG = sum(PROD_QTY)/uniqueN(LYLTY_CARD_NBR)), .(LIFESTAGE, PREMIUM_CUSTOMER)][ord

ggplot(data = avg_units, aes(weight = AVG, x = LIFESTAGE, fill = PREMIUM_CUSTOMER)) + geom_bar(position
labs(x = "Lifestage", y = "Avg units per transaction", title = "Units per customer") + theme(axis.text.
```
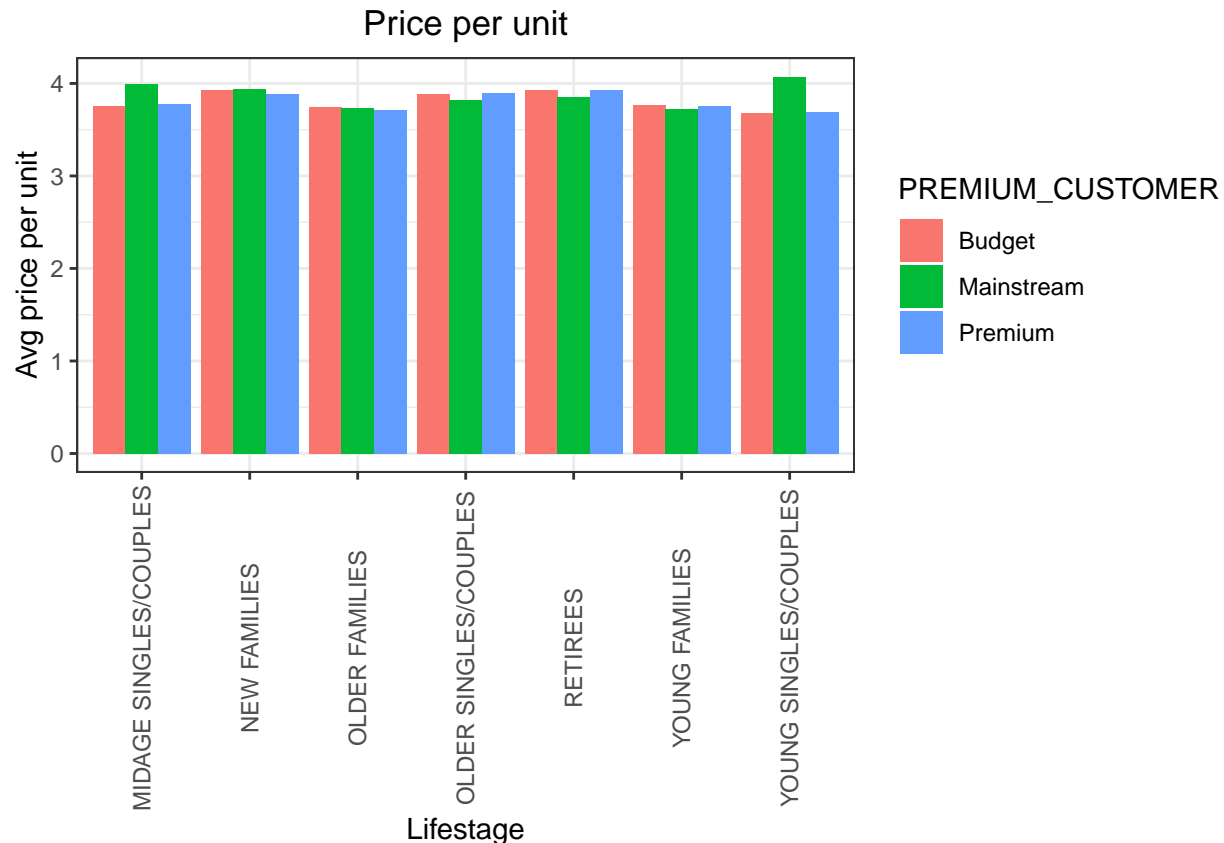
## Units per customer



Young families and old families have generally bought more chips in comparision with the midage and retirees

Lets investigate the average price per unit chip bought by each family

First compute average price per unit chips i.e total_sales/Prod_qty

```r
# Average price per unit by LIFESTAGE and PREMIUM_CUSTOMER
avg_price <- data[, .(AVG = sum(TOT_SALES)/sum(PROD_QTY)), .(LIFESTAGE, PREMIUM_CUSTOMER)][order(-AVG)]
#### Create plot
ggplot(data = avg_price, aes(weight = AVG, x = LIFESTAGE, fill = PREMIUM_CUSTOMER)) + geom_bar(position
```

## Price per unit



from the plot it is clear that mainstream- midage and young singles/couples are more willing to pay per packet of chips compared to budget and premium counterparts..

As the difference between the average price per unit is not same we can check this difference is stastically significant or not. . .

Performing independent t-test between mainstream vs premium and budget midage and young young single couples

```r
# young singles and couples

pricePerUnit <- data[, price := TOT_SALES/PROD_QTY]
t.test(data[LIFESTAGE %in% c("YOUNG SINGLES/COUPLES", "MIDAGE SINGLES/COUPLES") & PREMIUM_CUSTOMER == "
, data[LIFESTAGE %in% c("YOUNG SINGLES/COUPLES", "MIDAGE SINGLES/COUPLES") & PREMIUM_CUSTOMER != "Mains
, alternative = "greater")
```

```
##
##  Welch Two Sample t-test
##
## data:  data[LIFESTAGE %in% c("YOUNG SINGLES/COUPLES", "MIDAGE SINGLES/COUPLES") & PREMIUM_CUSTOMER ==
## t = 37.624, df = 54791, p-value < 0.00000000000000022
## 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
```

The t-test results in a p-value of 2.2e-16 , 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
segment1 <- data[LIFESTAGE == "YOUNG SINGLES/COUPLES" & PREMIUM_CUSTOMER == "Mainstream",]
other <- data[!(LIFESTAGE == "YOUNG SINGLES/COUPLES" & PREMIUM_CUSTOMER == "Mainstream"),]

## Brand affinity compared to the rest of the population
quantity_segment1 <- segment1[, sum(PROD_QTY)]

quantity_other <- other[, sum(PROD_QTY)]

quantity_segment1_by_brand <- segment1[, .(targetSegment = sum(PROD_QTY)/quantity_segment1), by = BRAND]

quantity_other_by_brand <- other[, .(other = sum(PROD_QTY)/quantity_other), by = BRAND]

brand_proportions <- merge(quantity_segment1_by_brand, quantity_other_by_brand)[, affinityToBrand := ta

brand_proportions[order(-affinityToBrand)]
```

```
##           BRAND targetSegment          other affinityToBrand
## 1:     TYRRELLS    0.031552795 0.025692464       1.2280953
## 2:     TWISTIES    0.046183575 0.037876520       1.2193194
## 3:      DORITOS    0.122760524 0.101074684       1.2145526
## 4:       KETTLE    0.197984817 0.165553442       1.1958967
## 5:      TOSTITOS    0.045410628 0.037977861       1.1957131
## 6:     PRINGLES    0.119420290 0.100634769       1.1866703
## 7:         COBS    0.044637681 0.039048861       1.1431238
## 8:     INFUZIONS    0.064679089 0.057064679       1.1334347
## 9:        THINS    0.060372671 0.056986370       1.0594230
## 10:     GRNWVES    0.032712215 0.031187957       1.0488733
## 11:    CHEEZELS    0.017971014 0.018646902       0.9637534
## 12:       SMITHS    0.096369910 0.124583692       0.7735355
## 13:       FRENCH    0.003947550 0.005758060       0.6855694
## 14:      CHEETOS    0.008033126 0.012066591       0.6657329
## 15:          RRD    0.043809524 0.067493678       0.6490908
## 16:      NATURAL    0.019599724 0.030853989       0.6352412
## 17:          CCS    0.011180124 0.018895650       0.5916771
## 18:     SUNBITES    0.006349206 0.012580210       0.5046980
## 19: WOOLWORTHS    0.024099379 0.049427188       0.4875733
## 20:       BURGER    0.002926156 0.006596434       0.4435967
```

We can see that : - Mainstream young singles/couples are 23% more likely to purchase Tyrrells chips compared to the rest of the population - Mainstream young singles/couples are 56% less likely to purchase Burger Rings compared to the rest of the population

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

```r
# Preferred pack size compared to the rest of the population
quantity_segment1_by_pack <- segment1[, .(targetSegment = sum(PROD_QTY)/quantity_segment1), by = PACK_S]
quantity_other_by_pack <- other[, .(other = sum(PROD_QTY)/quantity_other), by = PACK_SIZE]

pack_proportions <- merge(quantity_segment1_by_pack, quantity_other_by_pack)[, affinityToPack := targetS

pack_proportions[order(-affinityToPack)]
```

```
##     PACK_SIZE targetSegment       other affinityToPack
##  1:      270   0.031828847 0.025095929      1.2682873
##  2:      380   0.032160110 0.025584213      1.2570295
##  3:      330   0.061283644 0.050161917      1.2217166
##  4:      134   0.119420290 0.100634769      1.1866703
##  5:      110   0.106280193 0.089791190      1.1836372
##  6:      210   0.029123533 0.025121265      1.1593180
##  7:      135   0.014768806 0.013075403      1.1295106
##  8:      250   0.014354727 0.012780590      1.1231662
##  9:      170   0.080772947 0.080985964      0.9973697
## 10:      150   0.157598344 0.163420656      0.9643722
## 11:      175   0.254989648 0.270006956      0.9443818
## 12:      165   0.055652174 0.062267662      0.8937572
## 13:      190   0.007481021 0.012442016      0.6012708
## 14:      180   0.003588682 0.006066692      0.5915385
## 15:      160   0.006404417 0.012372920      0.5176157
## 16:       90   0.006349206 0.012580210      0.5046980
## 17:      125   0.003008972 0.006036750      0.4984423
## 18:      200   0.008971705 0.018656115      0.4808989
## 19:       70   0.003036577 0.006322350      0.4802924
## 20:      220   0.002926156 0.006596434      0.4435967
```

It looks like Mainstream young singles/couples are 27% more likely to purchase a 270g pack of chips compared to the rest of the population but let's dive into what brands sell this pack size.

## Conclusion

Let's recap what we've found! Sales have mainly been due to Budget - older families, Mainstream - young singles/couples, and Mainstream retirees shoppers. We found that the high spend in chips for mainstream young singles/couples and retirees is due to there being more of them than other buyers. Mainstream, midage and young singles and couples are also more likely to pay more per packet of chips. This is indicative of impulse buying behaviour. We've also found that Mainstream young singles and couples are 23% more likely to purchase Tyrrells chips compared to the rest of the population