

## Tasks:

- #### What is the count of Customers segmentation?
- #### What is the count of Customers attributes?
- #### What is the average total sales according to Customer segmentation?
- #### What is the average total sales according to Customers Life stage?
- #### What are the most favorite products brand and size?
- #### Which stores made the highest total sales?
- #### Which stores made the lowest total sales?
- #### What Customer type buys the most?
- #### At what date Sales are the highest?
- #### At what date transactions are the highest?

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sb
%matplotlib inline

In [2]: behav_data = pd.read_csv('QVI_purchase_behaviour.csv')
trans_data = pd.read_excel('QVI_transaction_data.xlsx')
```

```
Out [3]: behav_data.info()

Out[3]:
```

	LVLTY_CARD_NBR	LIFESTAGE	PREMIUM_CUSTOMER
0	1000	YOUNG SINGLES/COUPLES	Premium
1	1002	YOUNG SINGLES/COUPLES	Mainstream
2	1003	YOUNG FAMILIES	Budget
3	1004	OLDER SINGLES/COUPLES	Mainstream
4	1005	MIDAGE SINGLES/COUPLES	Mainstream
...	...	...	...
72632	2370651	MIDAGE SINGLES/COUPLES	Mainstream
72633	2370701	YOUNG FAMILIES	Mainstream
72634	2370751	YOUNG FAMILIES	Premium
72635	2370961	OLDER FAMILIES	Budget
72636	2373711	YOUNG SINGLES/COUPLES	Mainstream

72637 rows x 3 columns

```
In [4]: behav_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 72637 entries, 0 to 72636
Data columns (total 3 columns):
# Column Non-Null Count Dtype
---
0 LVLTY_CARD_NBR 72637 non-null int64
1 LIFESTAGE 72637 non-null object
2 PREMIUM_CUSTOMER 72637 non-null object
dtypes: int64(1), object(2)
memory usage: 1.7+ MB

In [5]: behav_data.LIFESTAGE.value_counts()

Out[5]:
```

LIFESTAGE	count
RETIRES	14805
OLDER SINGLES/COUPLES	14609
YOUNG SINGLES/COUPLES	14441
OLDER FAMILIES	9793
YOUNG FAMILIES	9178
MIDAGE SINGLES/COUPLES	9275
NEW FAMILIES	2549

Name: LIFESTAGE, dtype: int64

```
In [6]: behav_data.PREMIUM_CUSTOMER.value_counts()

Out[6]:
```

PREMIUM_CUSTOMER	count
Mainstream	29245
Budget	24470
Premium	18922

Name: PREMIUM\_CUSTOMER, dtype: int64

```
In [7]: behav_data.duplicated().unique()

Out[7]: array([False])

purchase_behaviour data is clean
```

```
In [8]: trans_data
```

```
Out[8]:
```

	DATE	STORE_NBR	LVLTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
0	43390	1	1000	1	5	Natural Chip Compy SeaSalt175g	2	6.0
1	43599	1	1307	348	66	CCs Nacho Cheese 175g	3	6.3
2	43605	1	1343	383	61	Smiths Crinkle Cut Chips Chicken 175g	2	2.9
3	43329	2	2373	974	69	Smiths Chip Thinly S/Cream&Onion 175g	5	15.0
4	43330	2	2426	1038	108	Kettle Tortilla ChpsHry&Jlono Chili 150g	3	13.8
...	...	...	...	...	...	...	...	...
264831	43533	272	27239	270058	89	Kettle Sweet Chili And Sour Cream 175g	2	10.8
264832	43225	272	272759	270164	74	Tostitos Splash Of Lime 175g	1	4.4
264833	43410	272	272739	270187	51	Doritos Mexicana 170g	2	8.8
264834	43461	272	272730	270188	42	Doritos Corn Chip Mexican Jalapeno 150g	2	7.8
264834	43465	272	272780	270198	74	Tostitos Splash Of Lime 175g	2	8.8

264836 rows x 8 columns

```
In [9]: trans_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 264836 entries, 0 to 264835
Data columns (total 8 columns):
# Column Non-Null Count Dtype
---
0 DATE 264836 non-null int64
1 STORE_NBR 264836 non-null int64
2 LVLTY_CARD_NBR 264836 non-null int64
3 TXN_ID 264836 non-null int64
4 PROD_NBR 264836 non-null int64
5 PROD_NAME 264836 non-null object
6 PROD_QTY 264836 non-null float64
7 TOT_SALES 264836 non-null float64
dtypes: float64(1), int64(6), object(1)
memory usage: 16.2+ MB

In [10]: trans_data
```

```
Out[10]:
```

	DATE	STORE_NBR	LVLTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_QTY	TOT_SALES
count	264836.000000	264836.000000	2.648330e+05	2.648330e+05	264836.000000	264833.000000	264833.000000
mean	43464.036700	135.080011	1.355490e+05	1.351583e+05	56.583167	1.907309	7.304200
std	105.389282	76.784189	8.057990e+04	7813305e+04	32.826638	0.643654	3.082326
min	43282.000000	1.000000	1.000000e+03	1.000000e+00	1.000000	1.000000	1.500000
25%	43373.000000	70.000000	7.002100e+04	6.760150e+04	28.000000	2.000000	5.400000
50%	43464.036700	130.000000	1.303570e+05	1.351370e+05	56.000000	2.000000	7.400000
75%	43555.000000	203.000000	2.030940e+05	2.027002e+05	85.000000	2.000000	9.200000
max	43646.000000	272.000000	2.373711e+06	2.415841e+06	114.000000	200.000000	650.000000

```
In [11]: trans_data.duplicated().value_counts()

Out[11]:
```

	count
False	264835
True	1

dtype: int64

```
In [12]: trans_data[trans_data.duplicated(keep=False)]

Out[12]:
```

	DATE	STORE_NBR	LVLTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
128483	43374	107	107024	108462	45	Smiths Thinly Cut Roast Chicken 175g	2	6.0
128485	43374	107	107024	108462	45	Smiths Thinly Cut Roast Chicken 175g	2	6.0

```
In [13]: trans_data.drop_duplicates(inplace=True)
trans_data.duplicated().value_counts()

Out[13]:
```

	count
False	264835
dtype: int64	

```
In [14]: trans_data.PROD_NAME.unique()

Out[14]:
```

array(['Natural Chip Compy SeaSalt175g', 'CCs Nacho Cheese 175g', 'Smiths Crinkle Cut Chips Chicken 170g', 'Smiths Chip Thinly S/Cream&Onion 175g', 'Kettle Tortilla ChpsHry&Jlono Chili 150g', 'Old El Paso Salsa Dip Tomato Mild 300g', 'Smiths Crinkle Chips Salt & Vinegar 330g', 'Grain Waves Sour Cream Mexican Jalapeno 150g', 'Doritos Corn Chip Mexican Jalapeno 150g', 'Grain Waves Sour CreamChives 210g', 'Fryella Crisps 150g', 'Twisties Cheese 270g', 'WW Crinkle Cut Chicken 175g', 'Things Chips Lights Tanga 175g', 'CCs Original 175g', 'Burger Rings 220g', 'Woolworths Cheese Rings 134g', 'Doritos Corn Chip Southern Chicken 150g', 'Cheesela Cheese Box 135g', 'Smiths Crinkle Infused Cn Crnchers Tanga Gnomole 110g', 'Kettle Sea Salt And Vinegar 175g', 'Infra Cn Crnchers Tanga Gnomole 110g', 'Red Rock Deli Thai Curry Original 175g', 'Kettle Original 175g', 'Pringles Stchn FriedChicken 134g', 'Pringles Sweet&Spicy BBQ 134g', 'Red Rock Deli SR Salsa & Mustila 150g', 'Things Chips Original salt 175g', 'Red Rock Deli Spn Salt & Truffle 150g', 'Smiths Crinkle Cut Original 175g', 'Kettle Chilli 175g', 'Doritos Mexicana 170g', 'Smiths Crinkle Cut French Onion 150g', 'Natural ChipCo Honey Soy Chkn175g', 'Dorito Corn Chip Supreme 380g', 'Twisties Chicken270g', 'Smiths Thinly Cut Roast Chicken 175g', 'Smiths Crinkle Cut Tomato Salsa 150g', 'Kettle Mozzarella Basil & Pesto 175g', 'Infusions Thai SweetChili Potato&On 110g', 'Kettle Sensations Caramel&T & Fig 150g', 'Smith Crinkle Cut Mac W Cheese 150g', 'Kettle Honey Soy Chicken 175g', 'Things Chips Seasonedchicken 175g', 'Smiths Crinkle Chips Salt & Vinegar 170g', 'Infusions BBQ Rib Frwn Crckers 110g', 'GrnWee Plus Btcoot & Chilli Jan 180g', 'Tyrells Crisps Lightly Salted 165g', 'Kettle Sweet Chilli And Sour Cream 175g', 'Doritos Salsa Medium 300g', 'Kettle 135g Gwt Pot Sea Salt', 'Pringles SourCream Onion 134g', 'Doritos Corn Chips Original 170g', 'Twisties Cheese Burger 270g', 'Old El Paso Salsa Dip ChknY Tom Rt 300g', 'Coba Popd Swt/Chili 435/Cream Chips 110g', 'Woolworths Mild Salsa 300g', 'Natural Chip Co Tmato Hrb&Spce 175g', 'Smiths Crinkle Cut Chips Original 175g', 'Coba Popd Sea Salt Chips 110g', 'Smiths Crinkle Cut Chips Ch&Onion170g', 'French Fries Potato Chips 110g', 'Old El Paso Salsa Dip Tomato Med 300g', 'Doritos Corn Chips Cheese Supreme 170g', 'Crisps 134g', 'Smiths Crinkle Cut Chips 110g', 'Kettle Chilli Coconut 150g', 'WW Original Corn Chips 200g', 'Things Potato Chips Hot & Spicy 175g', 'Coba Popd Sour Crm &Chives Chips 110g', 'Smiths Crinkle Chips Grnly Big Bag 380g', 'Doritos Corn Chips Nacho Cheese 170g', 'Kettle Sensations BBQ&Apple 150g', 'WW Dstyle Chip Sea Salt 200g', 'Pringles Chicken Salt Crisps 134g', 'WW Original Stacked Chips 160g', 'Smiths Chip Thinly S/Cut&Vinegr175g', 'Cheesels Cheese 330g', 'Things Chips Salt & Vinegar 175g', 'Tostitos Lightly Salted 175g', 'Things Chips Salt & Vinegar 175g', 'Smiths Crinkle Cut Chips Barbecue 170g', 'Cheetos Puffs 165g', 'Kettle Sweet Chilli & Sour Cream 165g', 'WW Crinkle Cut Original 175g', 'Tostitos Splash Of Lime 175g', 'Woolworths Medium Salsa 300g', 'Kettle Tortilla ChpsBtcoot&Ltoetta 150g', 'CCs Tasty Cheese 175g', 'Woolworths Cheese Rings 190g', 'Tostitos Smoked Chootille 175g', 'Pringles Barbecue 134g', 'WW Supreme Cheese Chicken 200g', 'Pringles Mystery Flavour 150g', 'Tyrells Crisps Ched & Chives 165g', 'Cheetos Ch & Bacon Balls 190g', 'Pringles Slit Vingar 134g', 'Coba Whlgrn Crisps Ched & Chives 165g', 'Infusions SourCream&Hbbs Veg Srce 110g', 'Kettle Tortilla ChpsFeta&Garlic 150g', 'Infusions Mango Chutny Papadum 70g', 'Doritos Salsa Medium 300g', 'Chickenorizi 150g', 'Kettle 165g', 'Red Honey Soy Chutney 150g', 'Sunbites Whlgrn Crisps Frch/Onion 90g', 'Red Salt & Vinegar 165g', 'Doritos Cheese Supreme 330g', 'Smiths Crinkle Cut Snags&Sauce 150g', 'WW Sour Cream &Onion&stacked Chips 165g', 'Natural ChipCo Sea Salt & Vinegar 175g', 'Red Rock Deli Chkn&Garlic Aioli 150g', 'Red SR Slow Rat Pork Belly 150g', 'RR Pot Sea Salt 165g', 'Smiths Crinkle Cut Bolognese 150g', 'Doritos Salsa Mild 300g', 'dtype:object)

```
In [15]: binsize = 10
bins = np.arange(0,trans_data['TOT_SALES'].max()+binsize, binsize)
plt.figure(figsize=(8,5))
plt.hist(data=trans_data, x = 'TOT_SALES',bins = bins);
```

```
Out [16]: trans_data[trans_data.TOT_SALES > 100].head()

Out[16]:
```

	DATE	STORE_NBR	LVLTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
69762	43331	226	226000	226201	4	Dorito Corn Chip Supreme 380g	200	650.0
69763	43605	226	226000	226210	4	Dorito Corn Chip Supreme 380g	200	650.0

We have 2 data tables that will affect our analysis , we need to drop them

```
In [17]: trans_data = trans_data[trans_data.TOT_SALES < 100]
trans_data.describe()
```

```
Out[17]:
```

	DATE	STORE_NBR	LVLTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_QTY	TOT_SALES
count	264833.000000	264833.000000	2.648330e+05	2.648330e+05	264833.000000	264833.000000	264833.000000
mean	43464.036700	135.079529	1.355490e+05	1.351577e+05	56.583398	1.905812	7.293931
std	105.389061	76.784189	8.058000e+04	7813305e+04	32.826498	0.343437	2.527244
min	43282.000000	1.000000	1.000000e+03	1.000000e+00	1.000000	1.000000	1.500000
25%	43373.000000	70.000000	7.002100e+04	6.760000e+04	28.000000	2.000000	5.400000
50%	43464.000000	130.000000	1.303570e+05	1.351370e+04	56.000000	2.000000	7.400000
75%	43555.000000	203.000000	2.030940e+05	2.027000e+05	85.000000	2.000000	9.200000
max	43646.000000	272.000000	2.373711e+06	2.415841e+06	114.000000	5.000000	29.500000

Changing the data format to a more readable format

```
In [18]: base_date = pd.Timestamp('1901-01-01')
trans_data = trans_data.copy()
trans_data.loc[:,['DATE','TXN_ID']] = base_date + pd.DateOffset(date_offset)
trans_data.head()
```

```
Out[18]:
```

	DATE	STORE_NBR	LVLTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
0	2019-08-20	1	1000	1	5	Natural Chip Compy SeaSalt175g	2	6.0
1	2020-05-15	1	1307	348	66	CCs Nacho Cheese 175g	3	6.3
2	2020-05-21	1	1343	383	61	Smiths Crinkle Cut Chips Chicken 175g	2	2.9
3	2019-08-20	2	2373	974	69	Smiths Chip Thinly S/Cream&Onion 175g	5	15.0
4	2019-08-20	2	2426	1038	108	Kettle Tortilla ChpsHry&Jlono Chili 150g	3	13.8

```
In [19]: trans_data['PROD_SIZE(g)'] = trans_data['PROD_NAME'].str.extract(r'([0-9]{1,3})')
trans_data['PROD_SIZE(g)'].value_counts()
```

```
Out[19]:
```

PROD_SIZE(g)	count
175	66339
150	41311
134	25102
110	22387
170	19983
165	15297
300	15166
90	12405
380	6416
270	6285
210	6272
200	4473
135	3257
250	2195
90	3008
150	2969
160	2970
220	1564
70	1007
180	1468
125	1454
Name: PROD_SIZE(g), dtype: int64	

```
In [20]: trans_data['PROD_BRAND'] = trans_data['PROD_NAME'].str.split(' ').str[0]
trans_data['PROD_BRAND'].value_counts().sort_index()
```

```
Out[20]:
```

PROD_BRAND	count
Burger	1554
CCs	4551
Cheetos	2927
Cheesels	4603
Cobs	9693
Doritos	3193
French	24862
French	1418
Grain	1272
GrnWee	1468
Infusions	11057
Infra	1464
Kettle	41288
McC	1419
Natural	6030
Old El Paso	9324
Grain Waves	7740
Natural ChipCo	7469
Tyrells	6442
CCs	4551
Sunbites	3008
Cheetos	2927
Burger Rings	1564
French Fries	1418
Name: PROD_BRAND, dtype: int64	

```
In [21]: trans_data['PROD_BRAND'].replace({'Dorito': 'Doritos','Red':'Red Rock Deli','RRD':'Red Rock Deli','WW':'Woolworths','Sunbites','Smith':'Smiths','Infra':'Infusions','McC':'Natural','Natural','Natural Chipco','Old':'Old El Paso','Grain':'Grain Waves','Gr','French':'French Fries','Burger': 'Burger Rings'},inplace=True)
trans_data['PROD_BRAND'].value_counts()
```

```
Out[21]:
```

PROD_BRAND	count
Kettle	41288
Smiths	31822
Doritos	28145
Pringles	25102
Red Rock Deli	1779
Woolworths	14757
Infusions	14201
Things	14075
Cobs	9693
Tostitos	9471
Pringles	9471
Old El Paso	9324
Grain Waves	7740
Natural ChipCo	7469
CCs	4551
Sunbites	3008
Cheetos	2927
Burger Rings	1564
French Fries	1418
Name: PROD_BRAND, dtype: int64	

```
In [22]: trans_data['STORE_NBR'] = 'Store ' + trans_data['STORE_NBR'].astype('str')

Now the transactions data is clean enough
```

## Explanatory analysis

Let's merge the two data frames into just one

```
In [23]: new_df = pd.merge(behav_data,trans_data, on='LVLTY_CARD_NBR')
new_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 264833 entries, 0 to 264832
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 LVLTY_CARD_NBR 264833 non-null int64
1 LIFESTAGE 264833 non-null object
2 DATE 264833 non-null datetime64[ns]
3 STORE_NBR 264833 non-null object
4 TXN_ID 264833 non-null int64
5 PROD_NBR 264833 non-null int64
6 PROD_NAME 264833 non-null object
7 PROD_QTY 264833 non-null float64
8 TOT_SALES 264833 non-null float64
9 PROD_BRAND 264833 non-null object
10 DATE_TIME 264833 non-null object
11 PROD_SIZE(g) 264833 non-null object
dtypes: datetime64[ns](1), float64(1), int64(4), object(6)
memory usage: 26.3+ MB
```

```
Out [24]: new_df
```

	LVLTY_CARD_NBR	LIFESTAGE	PREMIUM_CUSTOMER	DATE	STORE_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY
0	1000	YOUNG SINGLES/COUPLES	Premium	2019-10-19	Store 1	1	5	Natural Chip Compy SeaSalt175g	
1	1002	YOUNG SINGLES/COUPLES	Mainstream	2019-08-18	Store 1	2	58	Red Rock Deli Chkn&Garlic Aioli 150g	
2	1003	YOUNG FAMILIES	Budget	2020-03-08	Store 1	3	52	Grain Waves Sour Cream&Chives 210g	
3	1003	YOUNG FAMILIES	Budget	2020-03-09	Store 1	4	106	Natural ChipCo Honey Soy Chkn175g	
4	1004	OLDER SINGLES/COUPLES	Mainstream	2019-11-04	Store 1	5	96	WW Original Stacked Chips 160g	
...	...	...	...	...	...	...	...	...	...
264828	2370701	YOUNG FAMILIES	Mainstream	2019-12-10	Store 88	240378	24	Grain Waves Sweet Chili 210g	
264829	2370751	YOUNG FAMILIES	Premium	2019-10-03	Store 88	240394	60	Kettle Tortilla ChpsFeta&Garlic 150g	
264830	2370961	OLDER FAMILIES	Budget	2019-10-26	Store 88	240480	70	Tyrells Crps Lightly Salted 165g	
264831	2370961	OLDER FAMILIES	Budget	2019-10-29	Store 88	240481	65	Salda Dip Crnky Tom H&S0g	
264832	2373711	YOUNG SINGLES/COUPLES	Mainstream	2019-12-16	Store 88	241815	16	Smiths Crinkle Cut Chips Salt & Vinegar 330g	

264833 rows x 10 columns

Let's see the count of different customer types

```
In [28]: unique_members_df = new_df.drop_duplicates(subset=['LVLTY_CARD_NBR'])
plt.style.use('seaborn')
base_color = sb.color_palette()[0]
sb.histplot(unique_members_df['PREMIUM_CUSTOMER'],color=base_color)
```

```
plt.xlabel('Customer Segmentation',fontsize=14)
plt.xticks(fontsize=12)
plt.title('Customers count based on Segmentation',fontsize=14);
```

```
Out [28]:
```

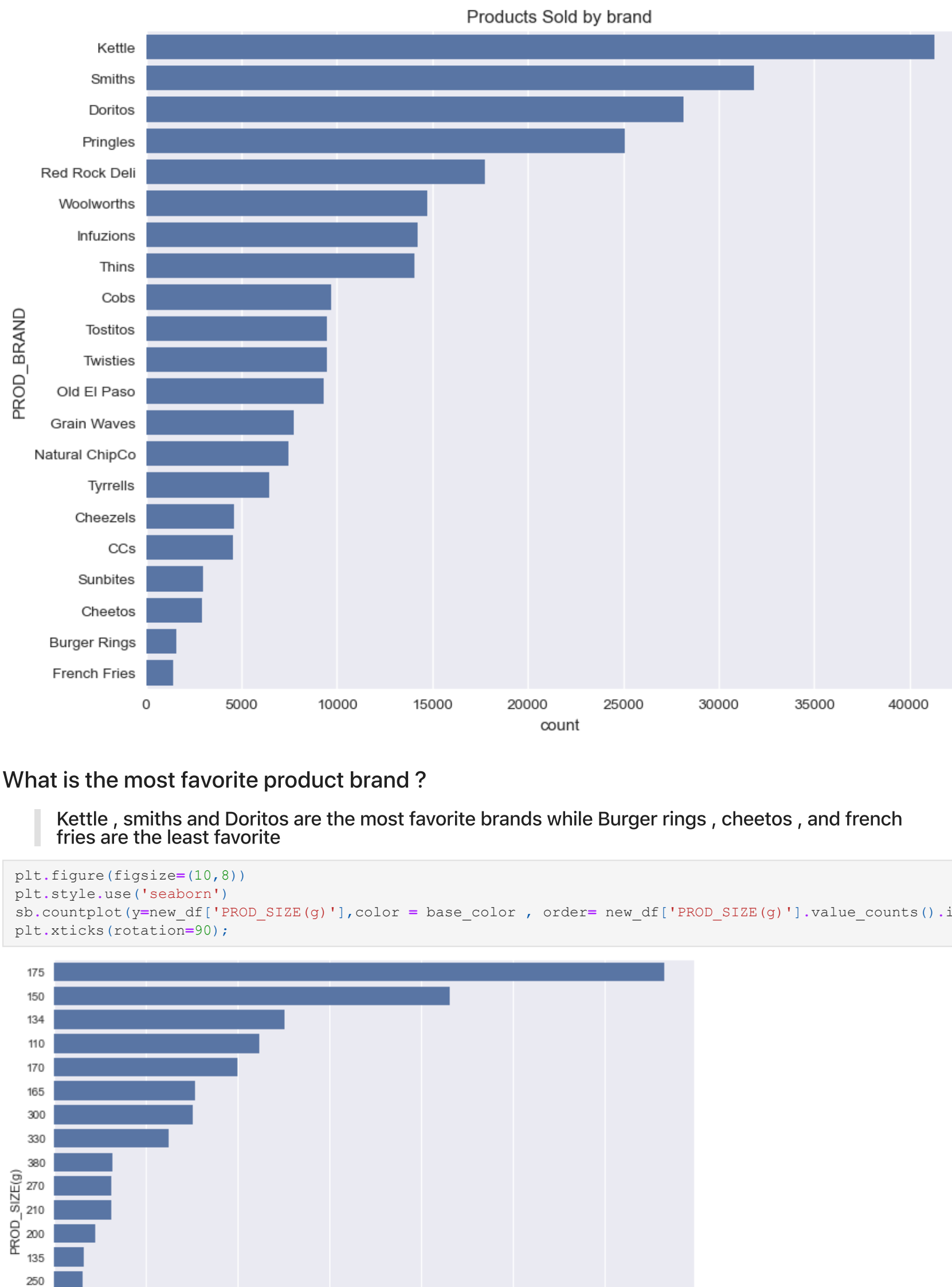
What is the count of Customers segmentation?

The mainstream customers are the most segmentation while the Premium customers are the lowest

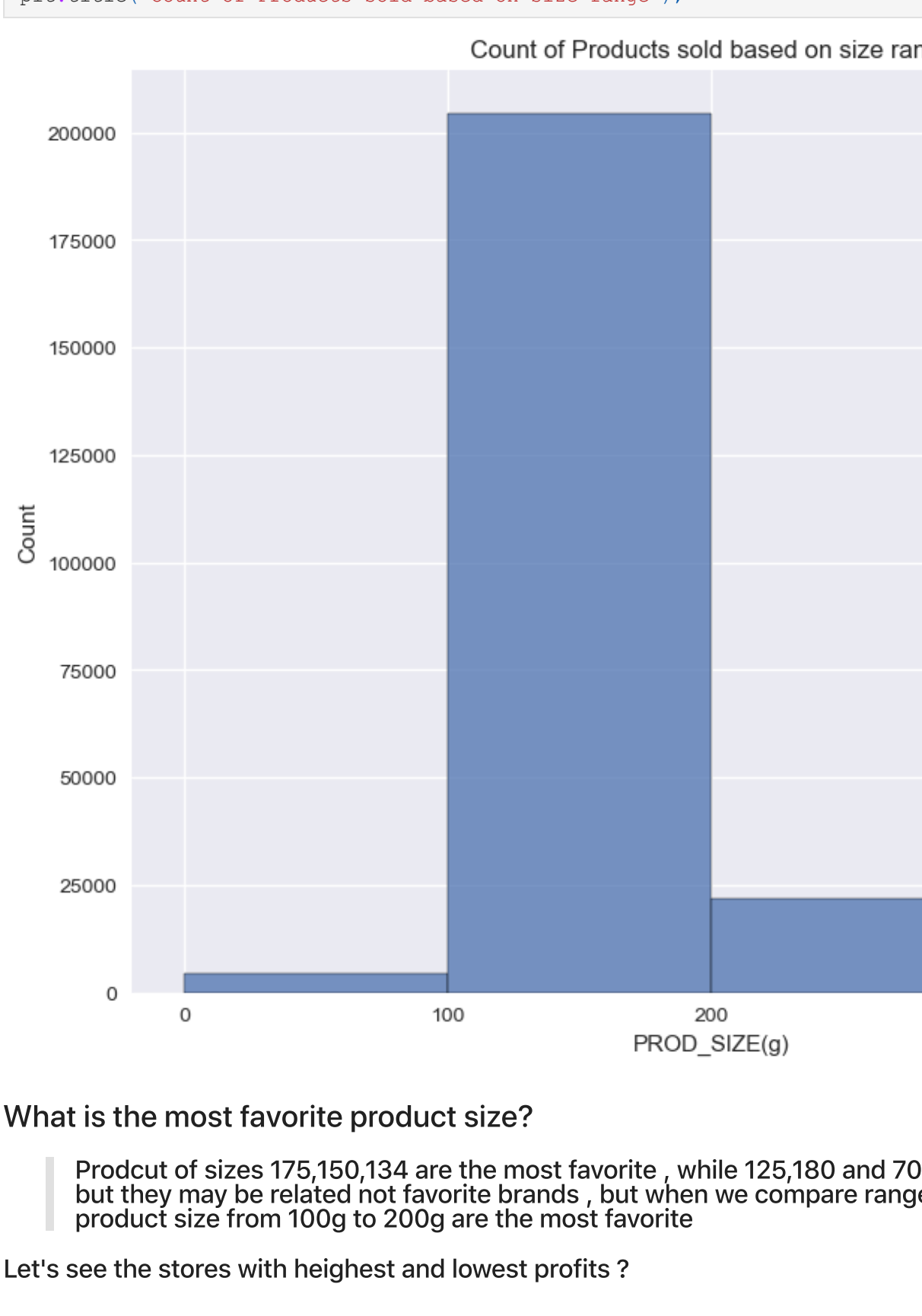
```
In [29]: plt.figure(figsize=(8,6))
plt.style.use('seaborn')

```





What is the most favorite product brand ?



What is the most favorite product size?

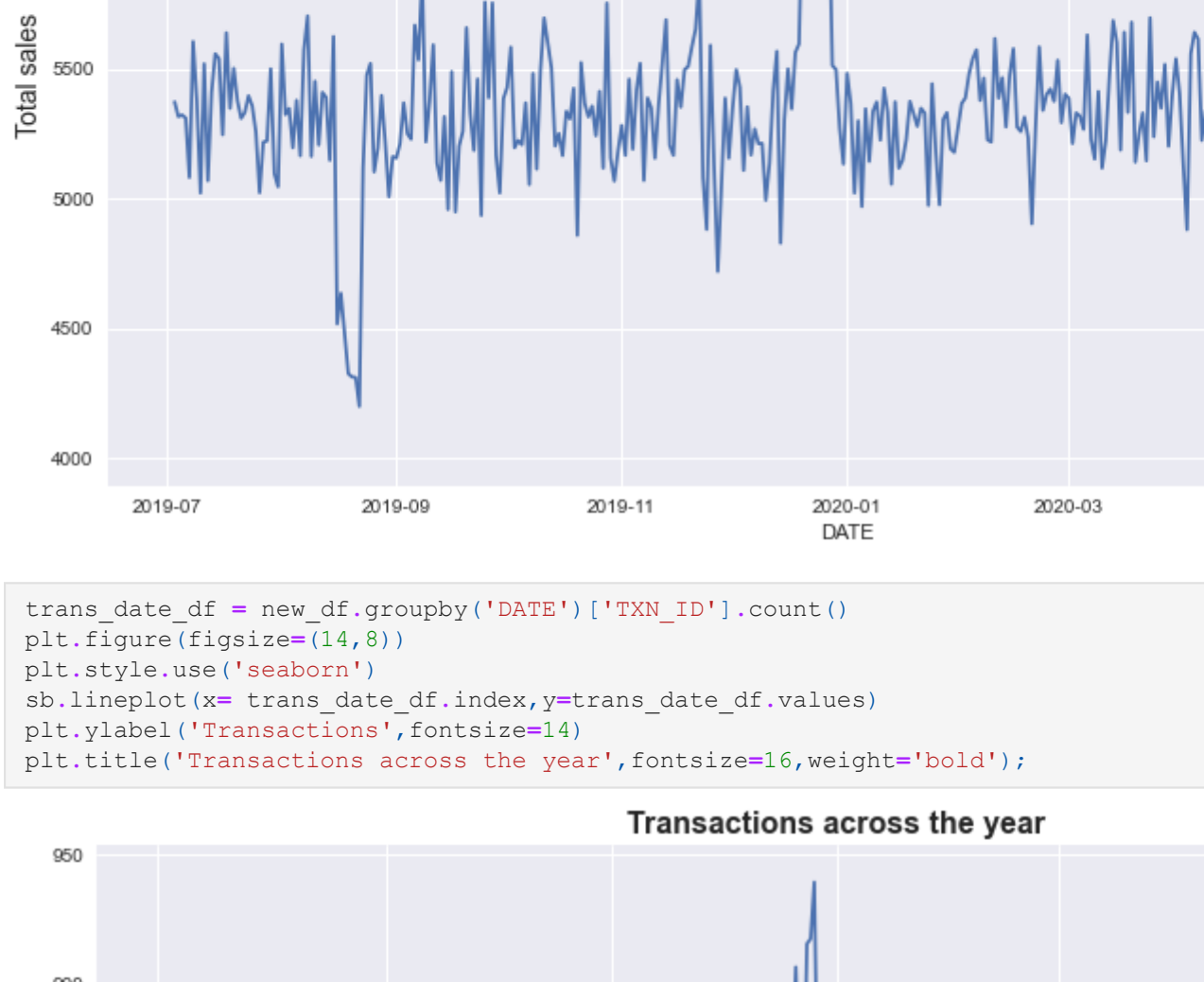
Product of sizes 175,150,134 are the most favorite , while 125,180 and 70g are the least favorite but they may be related not favorite brands , but when we compare ranges we will find that product size from 100g to 200g are the most favorite

Let's see the stores with heighest and lowest profits ?

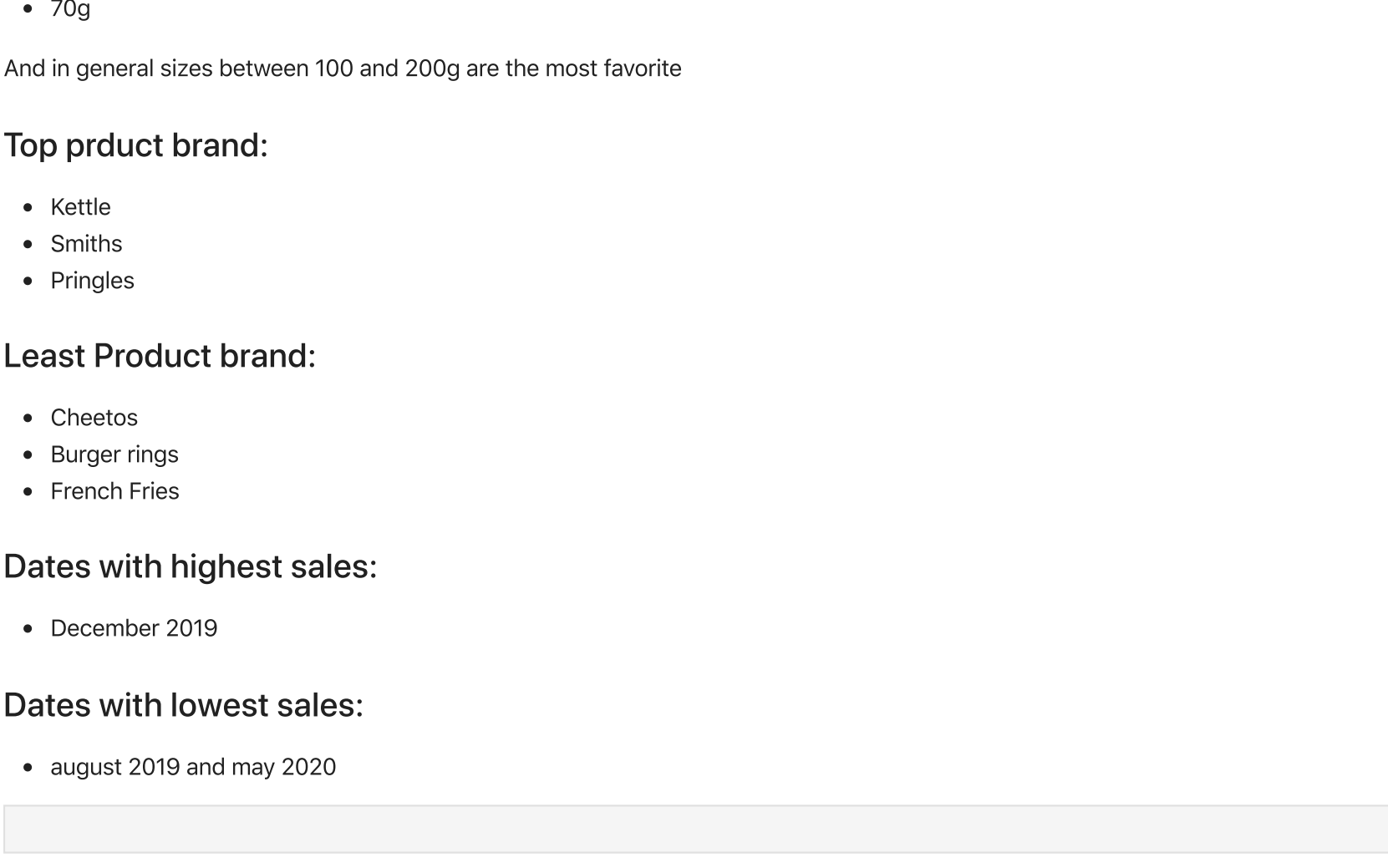
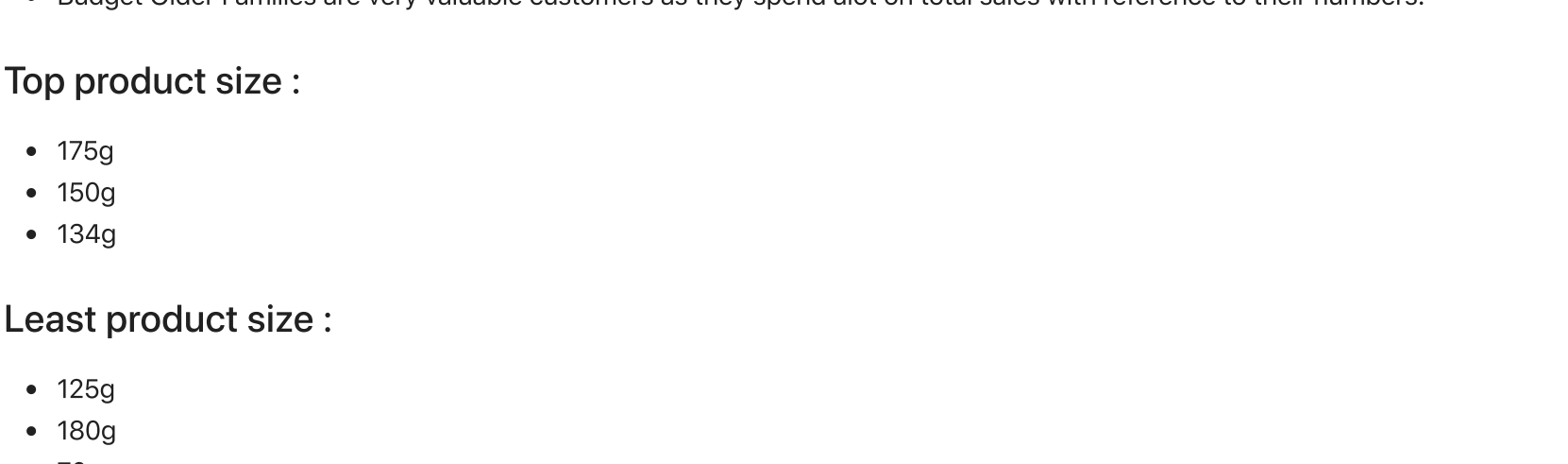
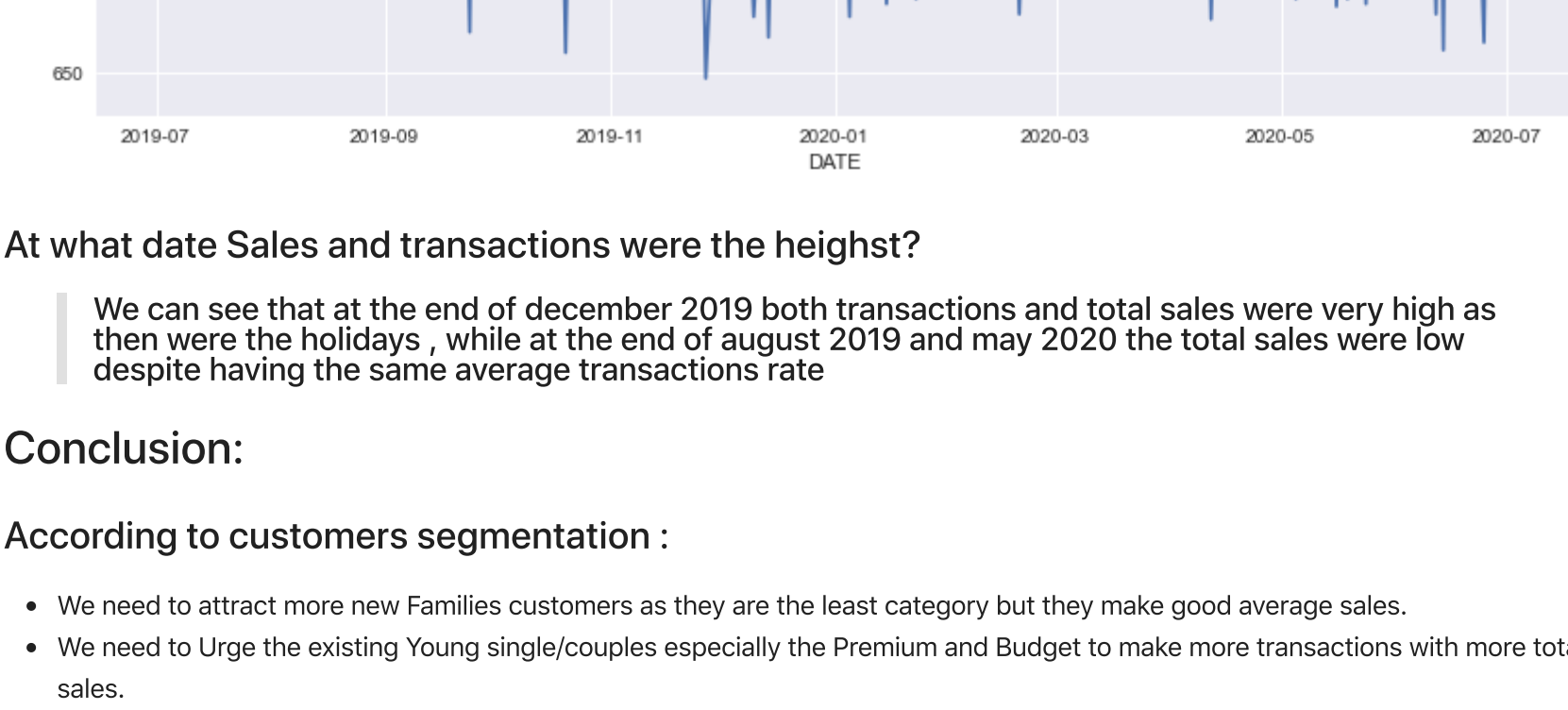
Which stores make the heighest total sales?



Which stores make the least total sales?



Now i want to analysis Sales and Transactions across the year to see when they were high or low?



At what date Sales and transactions were the heighest?

We can see that at the end of december 2019 both transactions and total sales were very high as then were the holidays , while at the end of august 2019 and may 2020 the total sales were low despite having the same average transactions rate

Conclusion:

According to customers segmentation :

- We need to attract more new Families customers as they are the least category but they make good average sales.
- We need to Urge the existing Young single/couples especially the Premium and Budget to make more transactions with more total sales.
- Budget Older Families are very valuable customers as they spend alot on total sales with reference to their numbers.

Top product size :

- 175g
- 150g
- 134g

Least product size :

- 125g
- 180g
- 70g

And in general sizes between 100 and 200g are the most favorite

Top prduct brand:

- Kettle
- Smiths
- Pringles

Least Product brand:

- Cheetos
- Burger rings
- French Fries

Dates with highest sales:

- December 2019

Dates with lowest sales:

- august 2019 and may 2020

