

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?
- #### 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 sns
import matplotlib inline

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

```
In [3]: behav_data

Out[3]:
```

	LYLVY_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	237051	MIDAGE SINGLES/COUPLES	Mainstream
72633	237071	YOUNG FAMILIES	Mainstream
72634	237051	YOUNG FAMILIES	Premium
72635	237061	OLDER FAMILIES	Budget
72636	237211	YOUNG SINGLES/COUPLES	Mainstream

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

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 72637 entries, 0 to 72636
Data columns (total 8 columns):
# Column Non-Null Count Dtype
---
0 LVLVY_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]: RETIRES 14805
OLDER SINGLES/COUPLES 14609
YOUNG SINGLES/COUPLES 14441
OLDER FAMILIES 9780
YOUNG FAMILIES 9178
MIDAGE SINGLES/COUPLES 7275
NEW FAMILIES 2549
Name: LIFESTAGE, dtype: int64
```

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

Out[6]: 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	LYLVY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
0	43390	1	1000	1	5	Natural Chip Compyr SeaSalt175g	2	6.0
1	43399	1	1307	348	66	CCs Nacho Cheese 175g	3	6.3
2	43205	1	1343	383	61	Smiths Crinkle Cut Chips Chicken 170g	2	2.9
3	43329	2	2428	1038	69	Smiths Chip Thiny S/Cream&Onion 175g	5	15.0
4	43330	2	2476	974	108	Kettle Tortilla ChpsHty&Jlpro Chll 150g	3	13.8
...
264631	43533	272	272319	270088	89	Kettle Sweet Chili And Sour Cream 175g	2	10.8
264632	43325	272	272356	270154	74	Toostitos Splash Of Lime 175g	2	4.4
264633	43410	272	272376	270187	51	Doritos Mexican 170g	1	8.8
264634	43461	272	272379	270188	42	Doritos Corn Chip Mexican Jalapeno 150g	2	7.8
264635	43365	272	272380	270189	74	Toostitos Splash Of Lime 175g	2	8.8

```
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 LVLVY_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.describe()

Out[10]:
```

	DATE	STORE_NBR	LYLVY_CARD_NBR	TXN_ID	PROD_NBR	PROD_QTY	TOT_SALES
count	264836.000000	264836.000000	2.648360e+05	2.648360e+05	264836.000000	264836.000000	264836.000000
mean	43464.036280	135.80811	1.355498e+05	1.351583e+05	56.583157	1.907309	7.304200
std	105.38982	76.78418	6.057998e+04	7.813303e+04	32.826638	0.643654	3.083226
min	43282.000000	1.000000	1.000000e+03	1.000000e+04	1.000000	1.000000	1.500000
25%	43373.000000	70.000000	7.002100e+04	6.760100e+04	28.000000	2.000000	5.400000
50%	43464.000000	130.000000	1.303575e+05	1.351375e+05	85.000000	2.000000	7.400000
75%	43555.000000	203.000000	2.030942e+05	2.027012e+05	66.000000	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]: False 264835
True 1
dtype: int64
```

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

Out[12]:
```

	DATE	STORE_NBR	LYLVY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
124843	43374	107	107024	108462	45	Smiths Thinly Cut Roast Chicken 175g	2	6.0
124845	43374	107	107024	108462	45	Smiths Thinly Cut Roast Chicken 175g	2	6.0

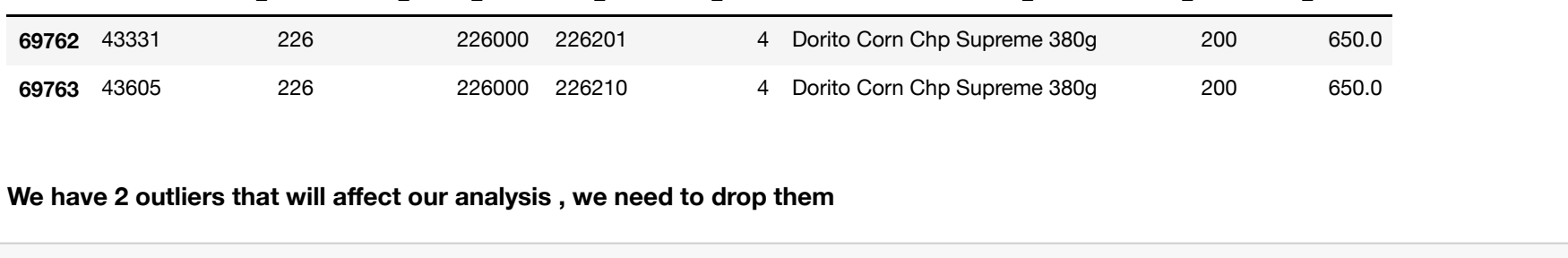
```
In [13]: trans_data.drop_duplicates(inplace=True)
trans_data.duplicated().value_counts()

Out[13]: False 264835
dtype: int64
```

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

Out[14]: array(['Natural Chip Compyr SeaSalt175g',
       'CCs Nacho Cheese 175g',
       'Smiths Crinkle Cut Chips Chicken 170g',
       'Smiths Chip Thiny S/Cream&Onion 175g',
       'Kettle Tortilla ChpsHty&Jlpro Chll 150g',
       'Old El Paso Salsa Dip Tomato Mild 300g',
       'Smiths Crinkle Chips Salt & Vinegar 330g',
       'Grain Waves Sweet Chili 210g',
       'Doritos Corn Chip Mexican Jalapeno 150g',
       'Grain Waves Sour Cream&Chives 210g',
       'Kettle Sensations Salsa Lime 150g',
       'Twisties Cheese 270g',
       'MW Crinkle Cut Chicken 175g',
       'Thins Chips Lightn Tangy 175g',
       'CCs Original 175g',
       'Burger Rings 220g',
       'MWCC Sour C Garden Chives 175g',
       'Doritos Corn Chip Southern Chicken 150g',
       'Cheezles Cheese Box 125g',
       'Smiths Crinkle Original 330g',
       'Infins Crn Cnchere Tangy Gcomale 110g',
       'Kettle Sea Salt And Vinegar 175g',
       'Smiths Chip Thiny Cut Original 175g',
       'Kettle Original 175g',
       'Red Rock Deli Thai ChilliLime 150g',
       'Fringles Ethn FriedChicken 134g',
       'Fringles Sweet&Spicy BBQ 134g',
       'Kettle Rock Deli Salsa & Mzzrila 150g',
       'Thins Chips Originl salted 175g',
       'Red Rock Deli Sp Salt & Truffle 150g',
       'Kettle Thiny Swt Chili 210g',
       'Doritos Mexicana 170g',
       'Smiths Crinkle Cut French OnionDip 150g',
       'Natural ChipCo Honey Soy Chokin150g',
       'Dorito Corn Chp Supreme 380g',
       'Twisties Chicken270g',
       'Smiths Thinly Cut Roast Chicken 175g',
       'Smiths Crinkle Cut Tomato Salsa 150g',
       'Kettle Mozzarella Basil & Pesto 170g',
       'Infuzions Thai SweetChili PotatoMix 110g',
       'Kettle Sensations Camembert & Fig 150g',
       'Smith Crinkle Cut Mac N Cheese 150g',
       'Kettle Honey Soy Chicken 175g',
       'Thins Chips SeasonedChicken 175g',
       'Smiths Crinkle Cut Salt & Vinegar 170g',
       'Kettle Sensations Salsa Lime 150g',
       'Infuzions BBQ Rib Prawn Crackers 110g',
       'GrnWves Plus Btroot & Chilli Jam 180g',
       'Tyrells Crisps Lightly Salted 165g',
       'Kettle Sweet Chili And Sour Cream 175g',
       'Doritos Salsa Medium 170g',
       'Kettle 135g Swt Pot Sea Salt',
       'Fringles SourCream Onion 134g',
       'Doritos Corn Chips Original 170g',
       'Twisties Cheese Burger 250g',
       'Old El Paso Salsa Dip Chnky Tom Ht300g',
       'Coba Popd Surt Czm 4Chives Chips 110g',
       'Kettle Tortilla Chl 4Sr/Cream Chips 110g',
       'Woolworths Mild Salsa 300g',
       'Natural Chip Co Tmato Hrb&Spce 175g',
       'Smiths Crinkle Cut Chips Original 170g',
       'Coba Popd Sea Salt Chips 110g',
       'Smiths Crinkle Cut Chips ChsOnion170g',
       'French Fries Potato Chips 175g',
       'Old El Paso Salsa Dip Tomato Med 300g',
       'Doritos Corn Chips Cheese Supreme 170g',
       'Fringles Original Crisps 134g',
       'MRD Chilli Coconut 150g',
       'MW Original Corn Chps Flavour 134g',
       'Thins Potate Chips Hot & Spicy 175g',
       'Coba Popd Sour Czm 4Chives Chips 110g',
       'Smiths Crinkle Cut Orgnl Big Bag 380g',
       'Doritos Corn Chips Nacho Cheese 170g',
       'Kettle Sensations BBQMaple 150g',
       'MW D/Style Chip Salt & Bag 380g',
       'Fringles Chicken Salt Crisps 134g',
       'MW Original Stacked Chips 160g',
       'Smiths Chip Thiny OutSalt/Vinegr175g',
       'Cheezles Cheese 330g',
       'Kettle Lightly Salted 175g',
       'Thins Chips Salt & Vinegar 175g',
       'Smiths Crinkle Cut Chips Barbecue 170g',
       'Cheetos Puffs 165g',
       'MRD Sweet Chili & Sour Cream 165g',
       'MW Crinkle Cut Original 175g',
       'Toostitos Splash Of Lime 175g',
       'Woolworths Medium Salsa 300g',
       'Kettle Tortilla ChpsBtroot&Ricotta 150g',
       'CCs Tasty Cheese 175g',
       'Woolworths Cheese Rings 190g',
       'Toostitos Smoked Chipotle 175g',
       'Fringles Barbeque 134g',
       'MW Supreme Cheese Corn Chips 200g',
       'Fringles Mystery Flavour 134g',
       'Tyrells Crisps Ched & Chives 165g',
       'Snbts Whgrn Crisps Cheddar&Metard 90g',
       'Cheetos Chs & Bacon Balls 190g',
       'Fringles Slc Vingar 134g',
       'Infuzions SourCream&Herbs Veg Strws 110g',
       'Kettle Tortilla ChpsFeta&Garlic 150g',
       'Infuzions Mango Chutny Papadums 70g',
       'MRD Steak & Chimichurri 150g',
       'MRD Honey Soy Chicken 165g',
       'Sunbites Whgrn Crisps Frch/Onion 90g',
       'MRD Salt & Vinegar 165g',
       'Doritos Cheese Supreme 330g',
       'Smiths Crinkle Cut Snags&Sauce 150g',
       'MW Sour Cream &OnionStacked Chips 160g',
       'MRD Lime & Pepper 165g',
       'Natural ChipCo Sea Salt & Vinegr 175g',
       'Red Rock Deli Chkn&Garlic Aioli 150g',
       'MRD SR Slow Roast Pork Belly 150g',
       'MRD Pc Sea Salt 165g',
       'Smith Crinkle Cut Bolognese 150g',
       'Doritos Salsa Mild 300g',
       dtype=object)
```

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



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

Out[16]:
```

	DATE	STORE_NBR	LYLVY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
69762	43301	226	226000	22601	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 outliers 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	LYLVY_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	43464.036570	135.079529	1.355489e+05	1.351577e+05	56.583598	1.905612	7.299351
std	105.389801	76.784189	6.058003e+04	7.813305e+04	32.826498	0.643437	3.082744
min	43282.000000	1.000000	1.000000e+03	1.000000e+04	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	1.303575e+05	1.351370e+05	85.000000	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 date format to a more readable format

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

Out[18]:
```

	DATE	STORE_NBR	LYLVY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
0	2019-05-15	1	1000	1	5	Natural Chip Compyr 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 170g	2	2.9
3	2019-08-19	2	2373	974	69	Smiths Chip Thiny S/Cream&Onion 175g	5	15.0
4	2019-08-20	2	2428	1038	108	Kettle Tortilla ChpsHty&Jlpro Chll 150g	3	13.8

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

Out[19]:
```

	LYLVY_CARD_NBR	LIFESTAGE	PREMIUM_CUSTOMER	DATE	STORE_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	PROD_I
0	1000	YOUNG SINGLES/COUPLES	Premium	2019-10-19	Store 1	1	5	Natural Chip Compyr SeaSalt175g	2	6.0
1	1002	YOUNG SINGLES/COUPLES	Mainstream	2019-09-18	Store 1	2	58	Red Rock Deli Chkn&Garlic Aioli 150g	2	6.0
2	1003	YOUNG FAMILIES	Budget	2020-03-08	Store 1	3	52	Grain Waves Sour Cream&Chives 210g	2	6.0
3	1003	YOUNG FAMILIES	Budget	2020-03-09	Store 1	4	106	Natural ChipCo Honey Soy Chokin175g	2	6.0
4	1004	OLDER SINGLES/COUPLES	Mainstream	2019-11-04	Store 1	5	96	MW Original Stacked Chips 160g	2	6.0
...
264828	2370701	YOUNG FAMILIES	Mainstream	2019-12-10	Store 88	240378	24	Grain Waves Sour Cream&Chives 210g	2	6.0
264829	2370751	YOUNG FAMILIES	Premium	2019-10-03	Store 88	240394	60	Kettle Tortilla ChpsFeta&Garlic 150g	2	6.0
264830	2370961	OLDER FAMILIES	Budget	2019-10-26	Store 88	240480	70	Tyrells Crisps Lightly Salted 165g	2	6.0
264831	2370961	OLDER FAMILIES	Budget	2019-10-29	Store 88	240481	65	Old El Paso Salsa Dip Chnky Tom Ht300g	2	6.0
264832	2373711	YOUNG SINGLES/COUPLES	Mainstream	2019-12-16	Store 88	241815	16	Smiths Crinkle Chips Salt & Vinegar 330g	2	6.0

264833 rows x 12 columns

Let's see the count of different customer types

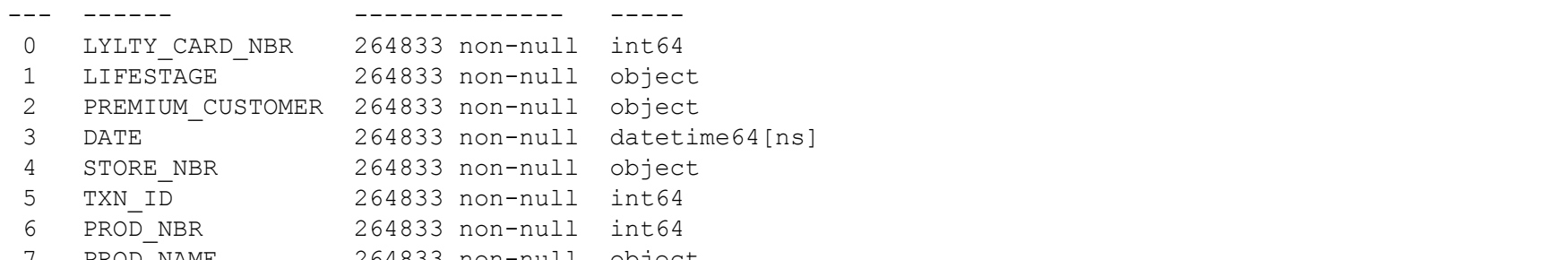
```
In [118]: unique_members_df = new_df.drop_duplicates(subset=['LYLVY_CARD_NBR'])
base_color = ab_color_palette()[0]
sb.histplot(unique_members_df['PREMIUM_CUSTOMER'])
plt.xlabel('Customer Segmentation',fontsize=14)
plt.xticks(fontsize=12)
plt.title('Customers count based on Segmentation',fontsize=14);
```



What is the count of Customers segmentation?

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

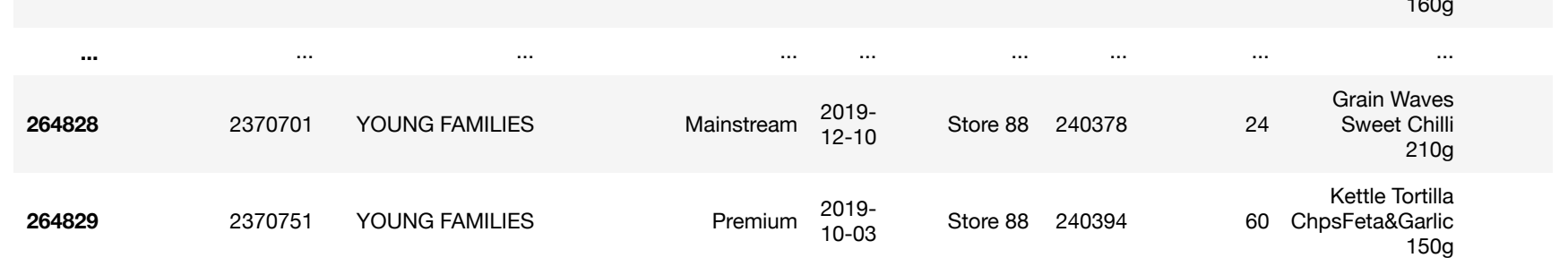
```
In [109]: plt.figure(figsize=(8,6))
plt.style.use('seaborn')
sb.countplot(y=new_df['LIFESTAGE'],order=new_df['LIFESTAGE'].value_counts().index, color = base_color)
plt.xlabel('Customers transactions based on Lifestage',fontsize=15);
plt.title('Customers transactions based on Lifestage',fontsize=15);
```



What is the average total sales according to Customer segmentation?

While the Older Families Customers are more than Older singles/couples but the last one made more transactions and more average total sales, while New Families are the least category with fewer transactions but they have high average total sales.

```
In [97]: plt.figure(figsize=(8,6))
base_color = ab_color_palette()[0]
lifest_data = new_df.groupby(y='LIFESTAGE')['TOT_SALES'].mean()
plt.style.use('seaborn')
sb.barplot(x=lifest_data.values,y = lifest_data.index, color = base_color)
plt.xlabel('Average total sales based on Lifestage',fontsize=14);
plt.title('Customers average total sales based on Lifestage',fontsize=15);
```

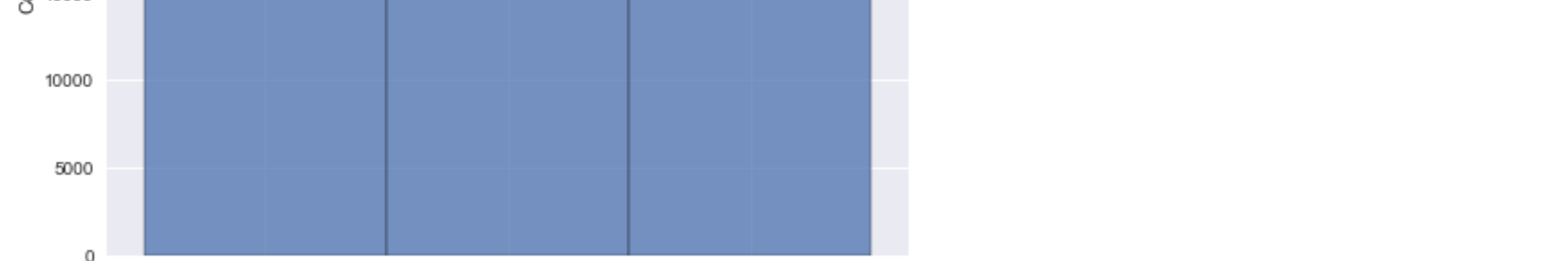


What is the average total sales according to Customers Life stage?

According to Lifestage Retires, Older and young single/couples are the most customers more than 14000 each, while New Families are the least about 2000

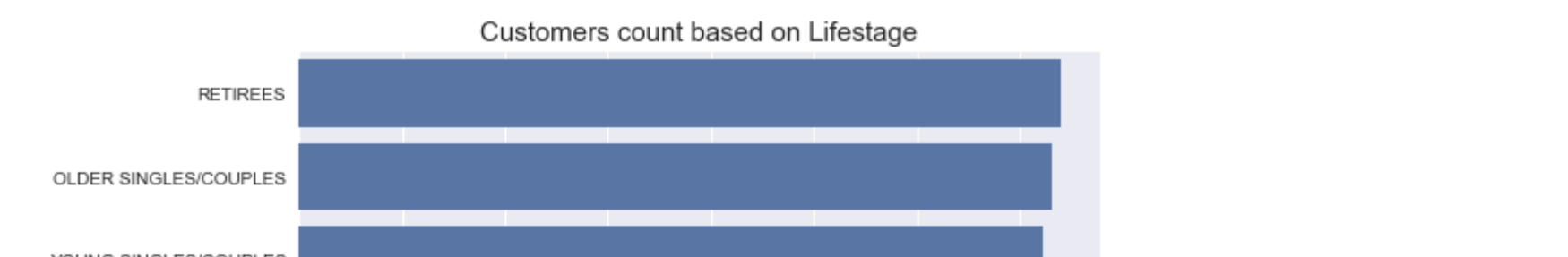
Now let's find out the average total sales according to the customer types

```
In [114]: plt.figure(figsize=(8,6))
base_color = ab_color_palette()[0]
customer_data = new_df.groupby(['PREMIUM_CUSTOMER'])['TOT_SALES'].mean()
plt.style.use('seaborn')
sb.barplot(x=customer_data.index,y = customer_data.values,color = base_color)
plt.xlabel('Average total sales based on Segmentation',fontsize=14);
plt.title('Customers average total sales based on Segmentation',fontsize=15);
```



what is the average total sales according to Customer segmentation?

```
In [110]: plt.figure(figsize=(8,6))
plt.style.use('seaborn')
sb.countplot(y=new_df['LIFESTAGE'],order=new_df['LIFESTAGE'].value_counts().index, color = base_color)
plt.xlabel('Customers transactions based on Lifestage',fontsize=15);
plt.title('Customers transactions based on Lifestage',fontsize=15);
```



```
In [97]: plt.figure(figsize=(8,6))
base_color = ab_color_palette()[0]
lifest_data = new_df.groupby(y='LIFESTAGE')['TOT_SALES'].mean()
plt.style.use('seaborn')
sb.barplot(x=lifest_data.values,y = lifest_data.index, color = base_color)
plt.xlabel('Average total sales based on Lifestage',fontsize=14);
plt.title('Customers average total sales based on Lifestage',fontsize=15);
```



what is the average total sales according to Customers Life stage?

However Retires Customers are more than Older singles/couples but the last one made more transactions and more average total sales, while New Families are the least category with fewer transactions but they have high average total sales.

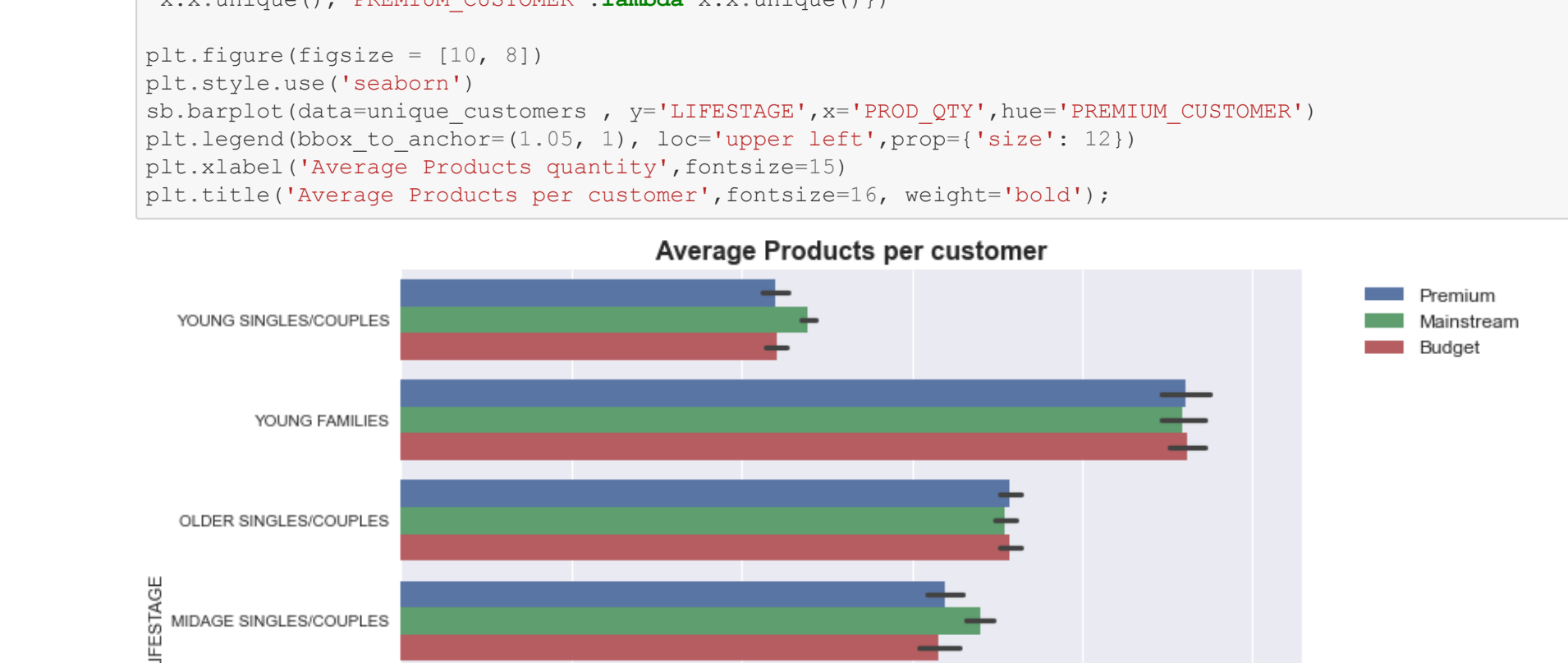
Now let's see the count of each customer lifestage and segmentation

```
In [95]: plt.figure(figsize=(12,8))
plt.style.use('seaborn')
sb.countplot(y=new_df['LIFESTAGE'],order=new_df['LIFESTAGE'].value_counts().index, hue = 'PREMIUM_CUSTOMER')
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left',prop={'size': 12})
plt.title('Customers count based on Lifestage',fontsize=15);
```

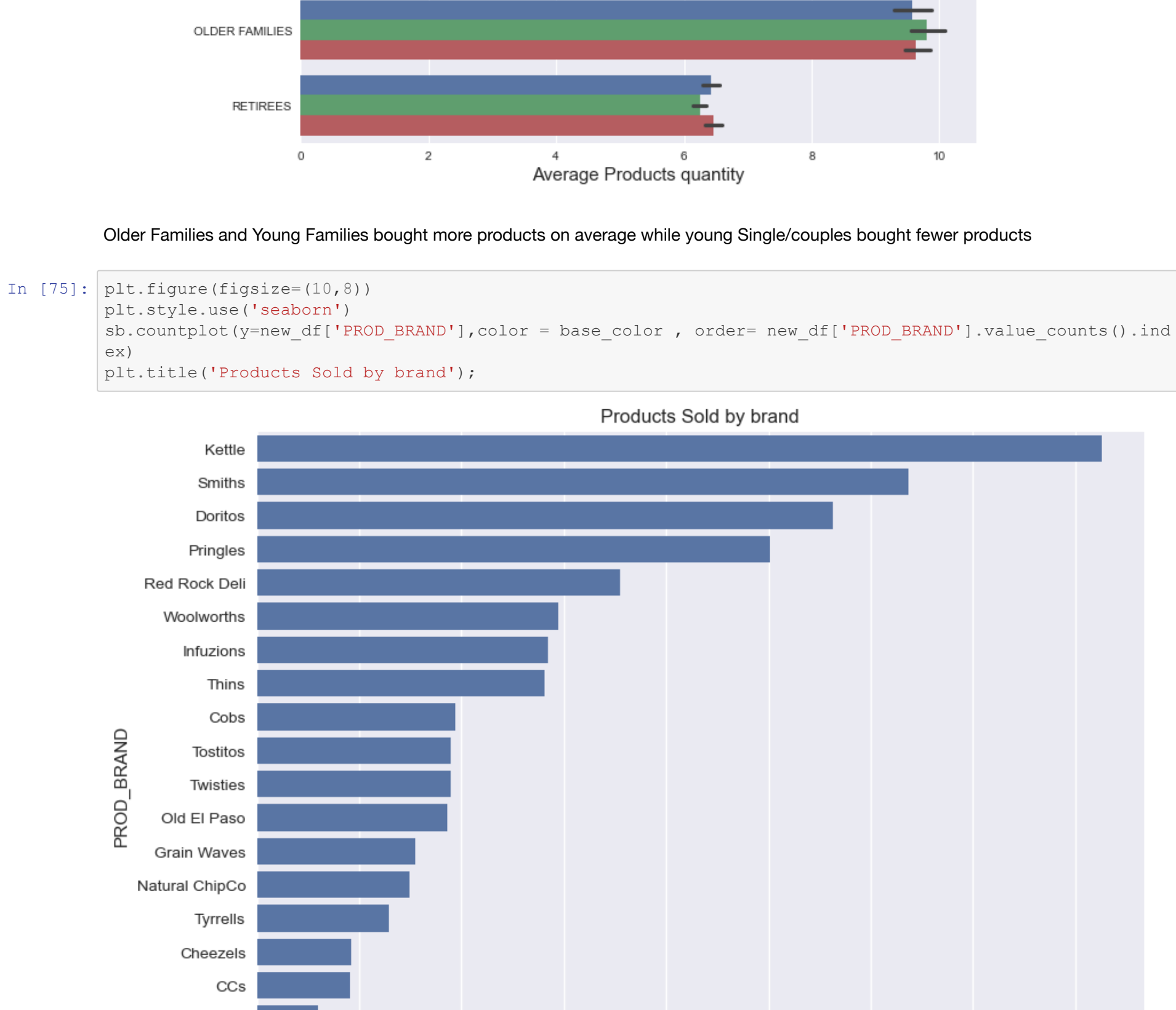

What is the average total sales per customer?

Young single/couples and Midage single/couples of mainstream segmentation make the highest average total sales while Young single/couples of budget and premium segmentation make the lowest average total sales

let's see How many products each Customer type bought

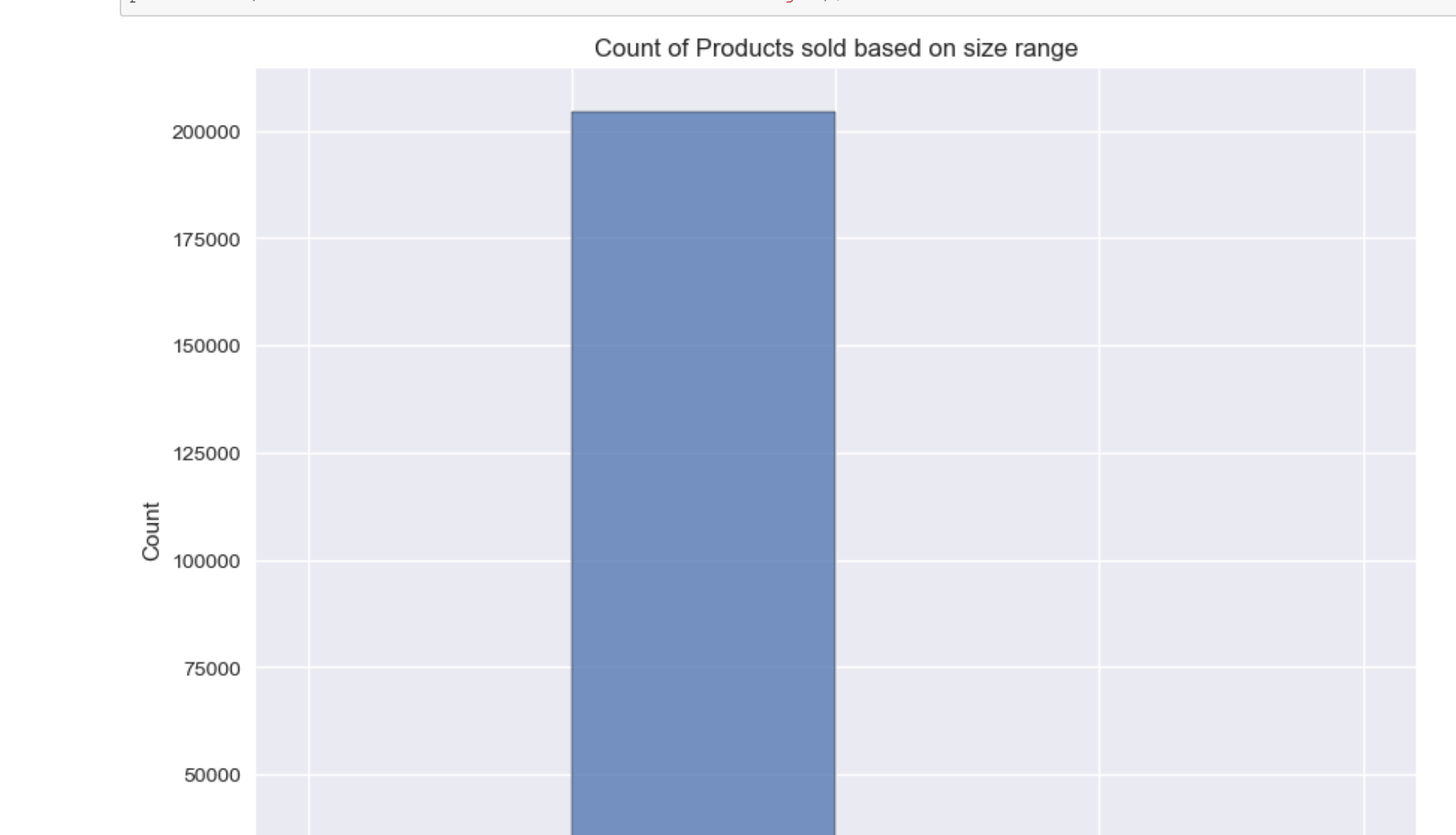
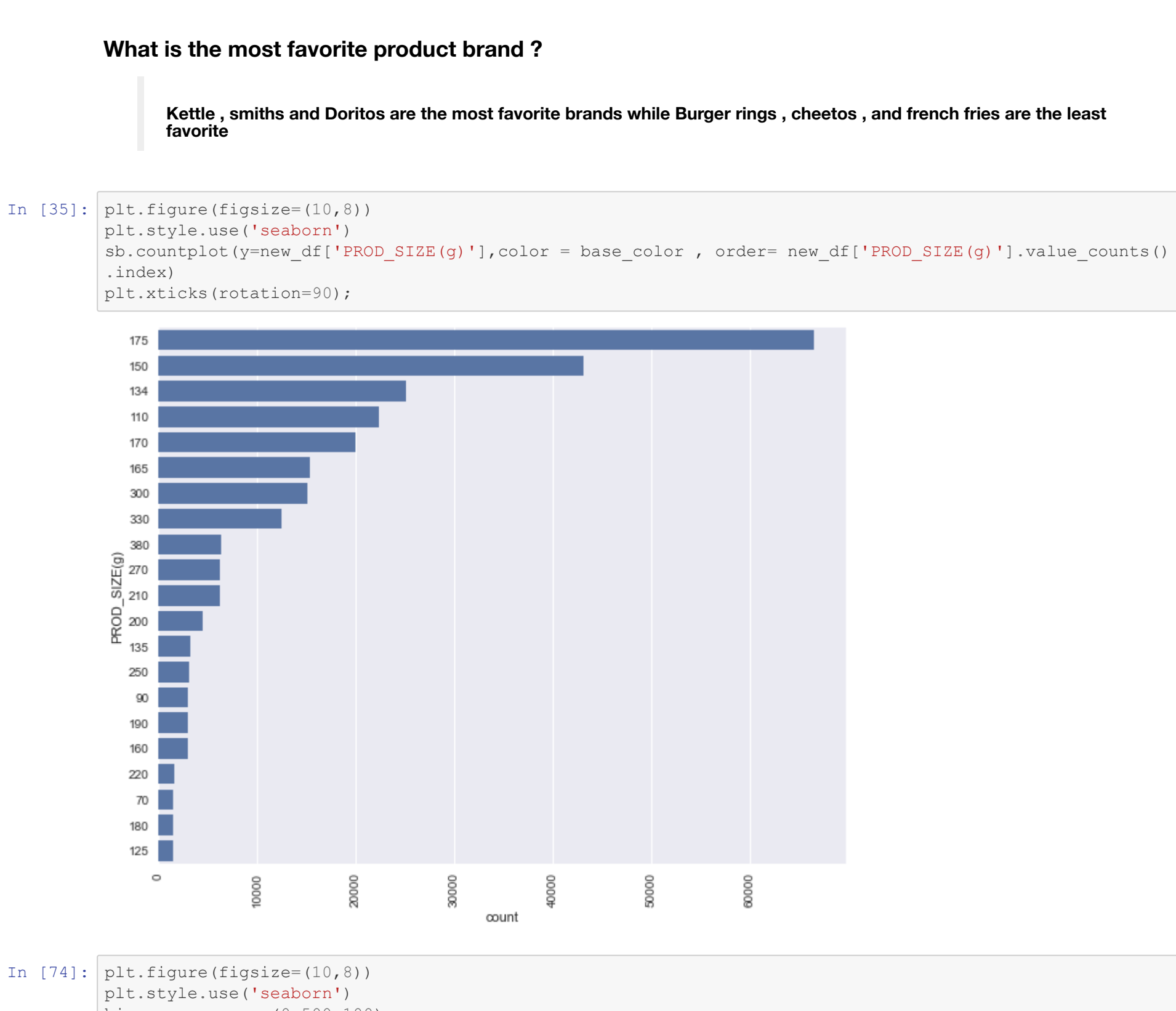


Older Families and Young Families bought more products on average while young Single/couples bought fewer products



What is the most favorite product brand ?

Kettle, smiths and Doritos are the most favorite brands while Burger rings , cheetos , and french fries are the least favorite

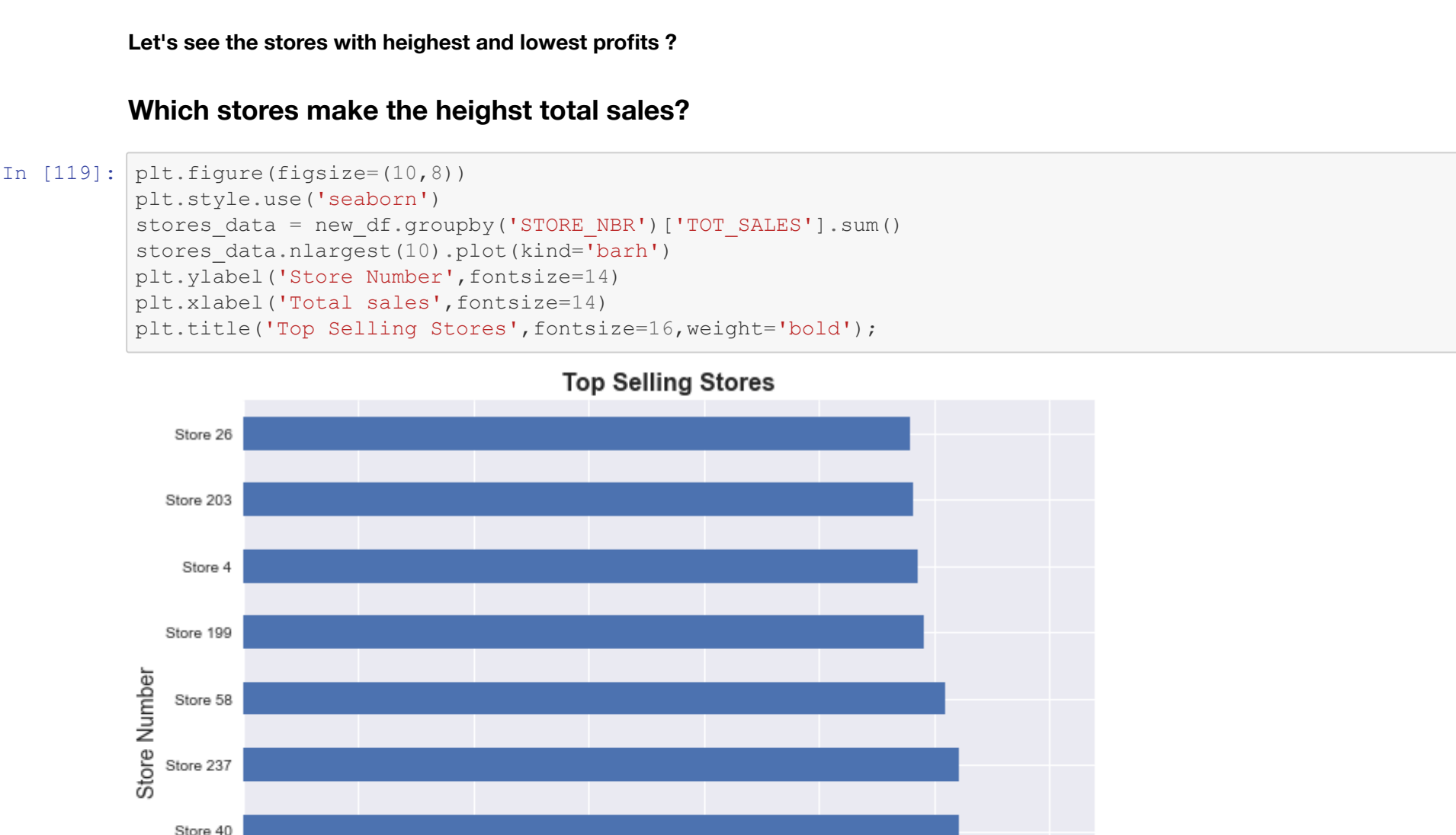


What is the most favorite product size?

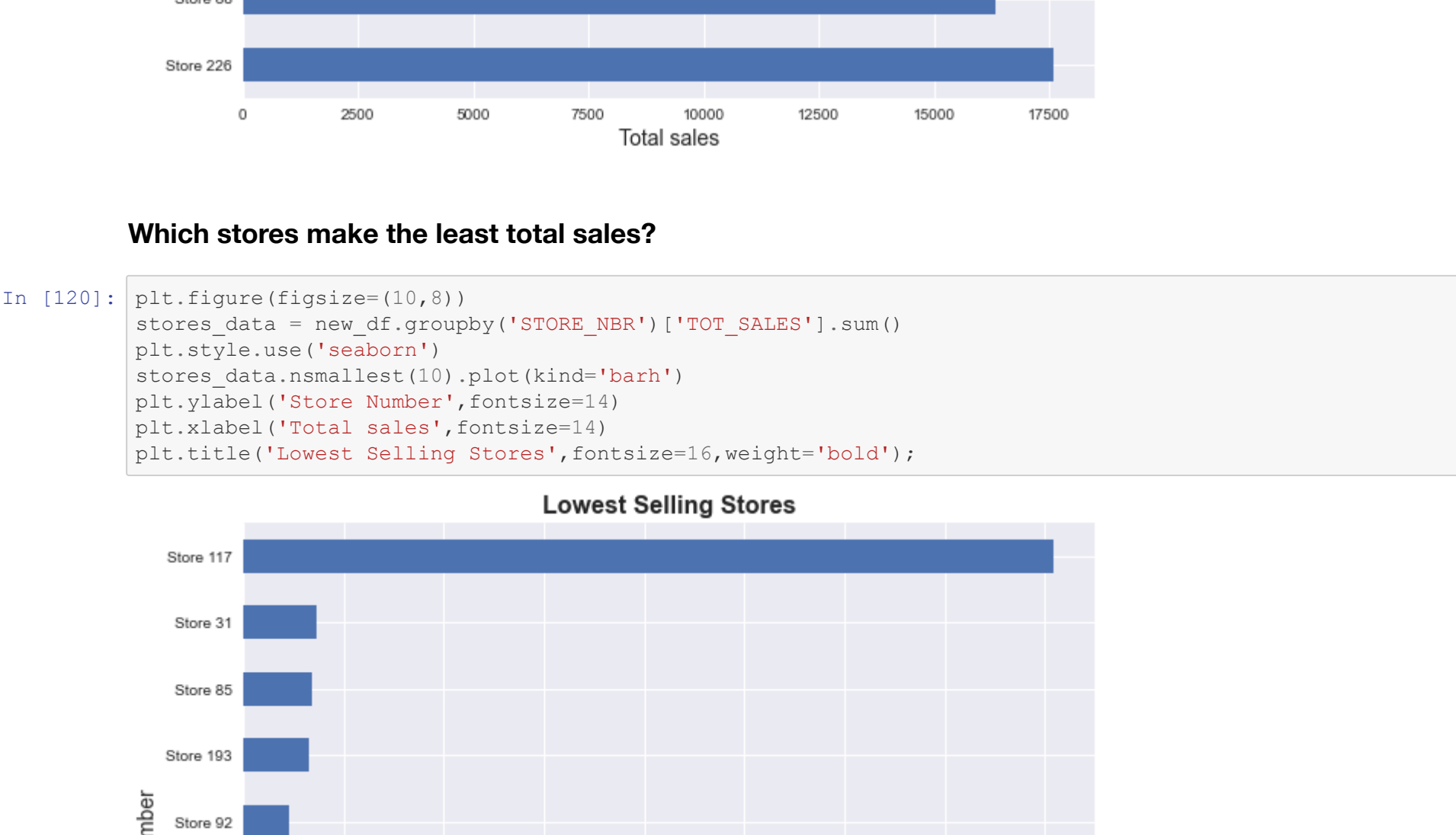
Product of sizes 175,150,134 are the most favorite , while 125,180 and 70g are the least favorite but they may be related not favorite brands , but when we compare ranges we will find that product size from 100g to 200g are the most favorite

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

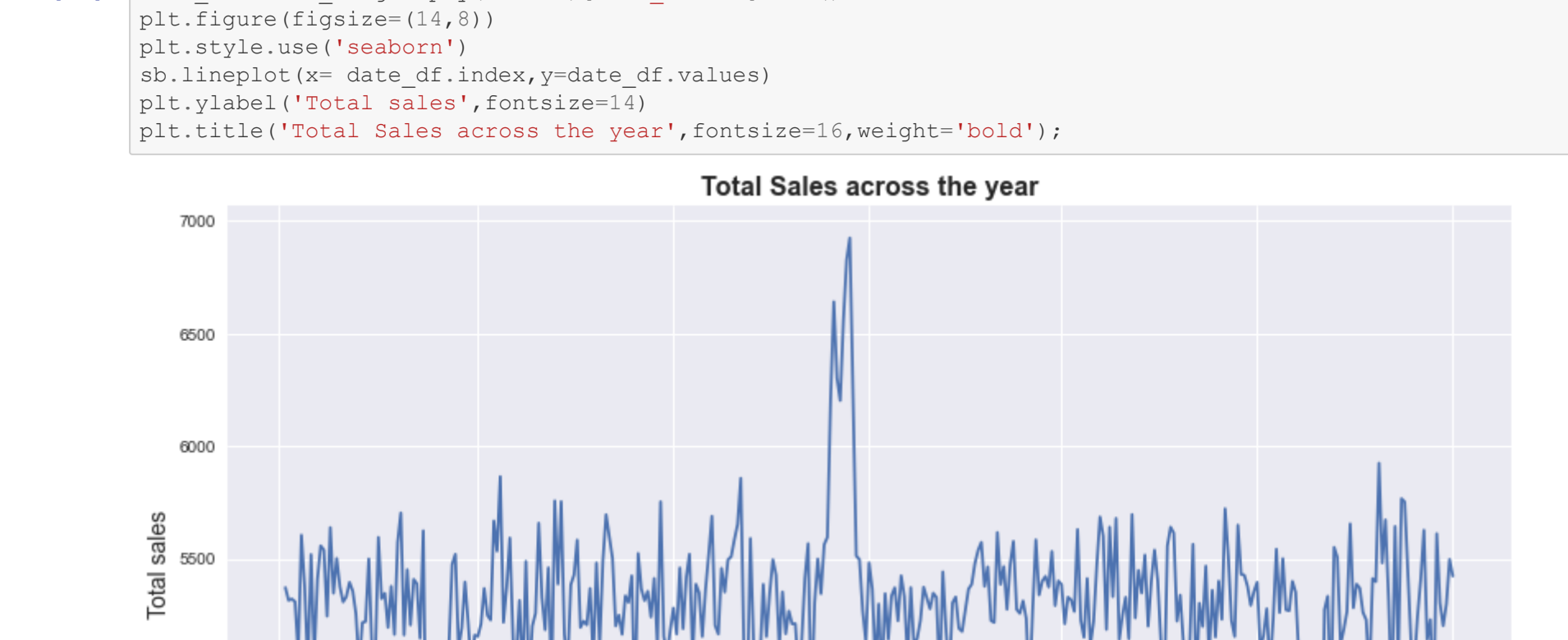
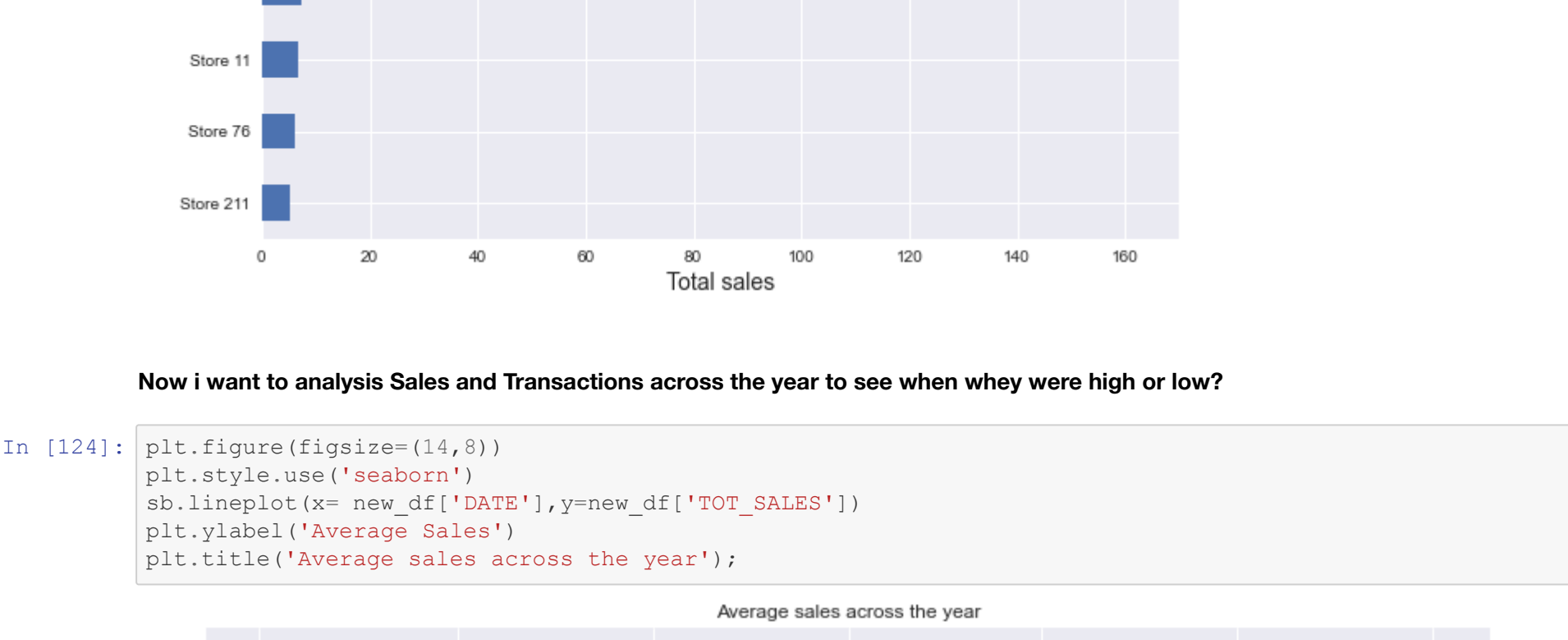
Which stores make the heighest total sales?



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

In [] :