

Chips Data Analysis

Introduction

The Quantum retail analytics team has been approached by the category manager of Chips, who wants to better understand the types of customers who purchase Chips and their purchasing behaviour within the region.

The insights from this analysis will feed into the supermarket's strategic plan for the chip category in the next half year.

In this analysis I will:

- Explore transaction data to look for inconsistencies, missing data, find outliers, correctly identify category items and numeric data across the different tables.
- Examine customer data to check for similar issues in customer data, look for nulls and merge the transaction data and customer data to prepare for data analysis
- Analyze data and customer segments by defining the metrics and looking at total sales, driver of sales, such as where the highest sales are coming from, etc. We will find interesting insights and trends in the data.
- Deep dive into customer segments by defining recommendations for the informed insights, determine which segments we should target, and if packet sizes and form and overall conclusion based on the analysis.

```
In [277]: #Import
import pandas as pd
import numpy as np

#Visualization
import seaborn as sns
import matplotlib.pyplot as plt

#Date
from datetime import datetime

#Warnings
import warnings
warnings.filterwarnings('ignore')

#Statistics
from scipy.stats import ttest_ind

#Regular Expression
import re

#Apriori
from mlxtend.frequent_patterns import apriori
from mlxtend.preprocessing import TransactionEncoder
from mlxtend.frequent_patterns import association_rules

#Others
from collections import Counter
```

```
In [2]: #Import datasets
transaction_data = pd.read_csv('QVI_transaction_data.csv')
customer_data = pd.read_csv('QVI_purchase_behaviour.csv')
```

Data Cleaning

Goals:

- Change date column to datetime
- Parse brand and pack sizes
- Remove special characters in PROD_NAME
- Remove sales brands from PROD_NAME
- Find outliers
- Replace brand names for similar brands

```
In [39]: customer_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 72637 entries, 0 to 72636
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   LIFESTAGE   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 [3]: #View dataset info
transaction_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 264835 entries, 0 to 264834
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   DATE        264835 non-null   datetime64[ns]
 1   STORE_NBR   264835 non-null   int64
 2   LIFESTAGE_NBR  264835 non-null   int64
 3   TXN_ID      264835 non-null   int64
 4   PROD_NBR    264835 non-null   int64
 5   PROD_NAME   264835 non-null   object
 6   PROD_QTY    264835 non-null   float64
 7   TOT_SALES   264835 non-null   float64
dtypes: float64(1), int64(6), object(1)
memory usage: 16.2+ MB
```

```
In [41]: transaction_data.head()
```

'Doritos Corn Chip Mexican Jalapeno 150g',	
'Grain Waves Sorp Cream&Chives 210g',	
'Kettle Sensations 150g',	'Siracha Lime 150g',
'Twisties Cheese 270g',	'WW Crinkle Cut Chicken 175g',
'Thin Chips Lightes 270g',	'CCs Original 175g',
'Burger Rings 220g',	'NOC Sorp Cream & Garden Chiles 175g',
'Doritos Corn Chip Southern Chichane 150g',	Original 330g',
'Smiths Crinkle Cut Chiles 125g',	'Smiths Crinkle
'Infense Crn Crnchers Tangy Gnoccole 110g',	
'Kettle Sea Salt And Vinegar 175g',	
'Kettle Chip Thiny Cut Original 175g',	'Kettle Original 175g',
'Red Rock Deli Thai ChilliLime 150g',	
'Fringles Strhn FriedChichane 134g',	'Fringles Sweet&Spicy BBQ 134g',
'Red Rock Deli Sr Salsa & Wazzila 150g',	
'Thin Chips Original asst 175g',	
'Red Rock Deli Sp Salt & Truffle 150g',	
'Smiths Thiny Swt Chills&Cream175g',	'Kettle Chilli 175g',
'Doritos Mexicana 175g',	
'Smiths Crinkle Cut Supreme OniOnDip 150g',	
'Natural ChipsCo Honey Soy Chkn175g',	
'Dorito Corn Chip French 380g',	'Twisties Chicken270g',

```
In [5]: #Converting serial date to datetime
date_changes = []
for i in transaction_data["DATE"]:
```

```
dt = datetime.fromordinal(datetime(1900,1,1).toordinal() + i - 1)
date_changes.append(dt)
transaction_data["DATE"] = date_changes
```

```
In [6]: transaction_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 264835 entries, 0 to 264834
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   DATE        264835 non-null   datetime64[ns]
 1   STORE_NBR   264835 non-null   int64
 2   LIFESTAGE_NBR  264835 non-null   int64
 3   TXN_ID      264835 non-null   int64
 4   PROD_NBR    264835 non-null   int64
 5   PROD_NAME   264835 non-null   object
 6   PROD_QTY    264835 non-null   float64
 7   TOT_SALES   264835 non-null   float64
dtypes: datetime64[ns](1), float64(1), int64(5), object(1)
memory usage: 16.2+ MB
```

```
In [7]: #Check that we are looking at the right products
transaction_data["PROD_NAME"].unique()
```

```
Out[7]: array(['Natural Chip Compy SeaSalt175g',
       'CCs Nacho Cheese 175g',
       'Smiths Crinkle Cut Chips Chicken 170g',
       'Smiths Chip Thiny S/Cream&Onion 175g',
       'Kettle Tortilla Chps&ny&lpno Chll 150g',
       'Old El Paso Salsa Dip Tomato Mild 300g',
       'Smiths Crinkle Chips Salt & Vinegar 330g',
       'Grain Waves Sweet Chilli 210g',
       'Doritos Corn Chip Mexican Jalapeno 150g',
       'Grain Waves Sour Cream&Chives 210g',
       'Kettle Sensations Siracha Lime 150g',
       'Twisties Cheese 270g', 'WW Crinkle Cut Chicken 175g',
       'Thins Chips Lightn Tangy 175g', 'Cce Original 175g',
       'Burger Rings 220g', 'NCC Sour Cream & Garden Chives 175g',
       'Doritos Corn Chip Southern Chicken 150g',
       'Cheezels Cheese Box 125g', 'Smiths Crinkle Original 330g',
       'Infins crn Cnchere Tangy Gamole 175g',
       'Kettle Sea Salt And Vinegar 175g',
       'Smiths Chip Thiny Cut Original 175g',
       'Red Rock Deli Thai Chilli&Lime 150g',
       'Pringles Stn FriedChicken 134g', 'Pringles Sweet&Spicy BBQ 134g',
       'Red Rock Deli SR Salsa & Mzzrila 150g',
       'Thins Chips Original salt&d 175g',
       'Red Rock Deli Sp Salt & Truffle 150g',
       'Smiths Thiny Swt Chl&d/Cream175g', 'Kettle Chilli 175g',
       'Doritos Mexicana 170g',
       'Smiths Crinkle Cut French Onion&Pip 150g',
       'Natural ChipCo Rony Tony Chok170g',
       'Dorito Corn Chip Supreme 380g', 'Twisties Chicken270g',
       'Smiths Thiny Cut Roast Chicken 175g',
       'Smiths Crinkle Cut Tomato Salsa 150g',
       'Kettle Mozzarella Basil & Pesto 175g',
       'Infuzions Thai Sweet&Chili Pot&toMx 110g',
       'Kettle Sensations Casm&e 175g',
       'Smith Crinkle Cut Mac N Cheese 150g',
       'Kettle Honey Soy Chicken 175g',
       'Thins Chips Seasoned&Chives 175g',
       'Smiths Crinkle Cut Salt & Vinegar 170g',
       'Infuzions BBQ Rib Prawn Crackers 110g',
       'GrnVes Plus Btroot & Chilli Jam 180g',
       'Pringles Stn FriedChicken 134g', 'Pringles Sweet&Spicy 175g',
       'Kettle Sweet Chilli And Sour Cream 175g',
       'Doritos Salsa Medium 300g', 'Kettle 135g Swt Pot Sea Salt',
       'Pringles SouCream Onion 134g',
       'Doritos Corn Chips Original 170g',
       'Twisties Cheese Burger 250g',
       'Old El Paso Salsa Dip Chnky Tom Ht300g',
       'Coba Popd Swt/Chilli Str&Cream&Pip 110g',
       'Woolworths Mild Salsa 110g',
       'Smiths Crinkle Cut Chips Chs&Onion170g',
       'French Fries Potato Chips 175g',
       'Old El Paso Salsa Dip Tomato Med 300g',
       'Doritos Corn Chips Cheese Supreme 170g',
       'Pringles Original Crisps 134g',
       'PRD Chilli Mango Coconut 150g',
       'WW Original Corn Chips 200g',
       'Thins Potato Chips Hot & Spicy 175g',
       'Coba Popd Sour Crm Chives Chips 110g',
       'Smiths Crinkle Cut Orgnl Big Bag 380g',
       'Doritos Corn Chips Nacho Cheese 170g',
       'Kettle Sensations BBQ&Maple 150g',
       'WW Style Chip Sea Salt 200g',
       'Pringles Chicken Salt Crisps 134g',
       'WW Original Stacked Chips 160g',
       'Smiths Chip Thiny Cut&SaltVineg175g', 'Cheezels Cheese 330g',
       'Thintos Lightly Salted 175g',
       'Thins Chips Salt & Vinegar 175g',
       'Smiths Crinkle Cut Chips Barbecue 170g', 'Cheetos Puffs 165g',
       'PRD Sweet Chilli & Sour Cream 165g',
       'WW Crinkle Cut Original 175g',
       'Twisties Splash Of Lime 175g', 'Woolworths Medium Salsa 300g',
       'Kettle Tortilla Chps&ny&lpno Chll 150g',
       'CCs Tasty Cheese 175g', 'Woolworths Cheese Rings 190g',
       'Twisties Smoked Chiptle 175g', 'Pringles Barbecue 134g',
       'WW Supreme Cheese Corn Chips 200g',
       'Pringles Mxy Flavour 134g',
       'Pyrellis Bacon Ched & Chives 160g',
       'Snts Whgrn Crisps Ched&Mtrd 90g',
       'Cheetos Chs & Bacon Btts 180g', 'Pringles Slit Vingar 134g',
       'Infuzions SourCream&Herbs Veg Strws 110g',
       'Kettle Tortilla Chps&ny&lpno 150g',
       'Pyrellis Mango Chutny Papadums 70g',
       'PRD Steak & ChimuChurri 150g',
       'PRD Honey Soy Chicken 160g',
       'Bumites Whgrn Crisps F&to170g',
       'PRD Salt & Vinegar 165g', 'Doritos Cheese Supreme 330g',
       'Smiths Crinkle Cut Snags&Sauce 150g',
       'WW Sour Cream &Onion&Stacked Chips 160g',
       'PRD Lime & Pepper 165g',
       'Natural ChipCo Sea Salt & Vineg 175g',
       'Red Rock Deli Chkn&Garlic Aioli 150g',
       'PRD SR Slow Roast Beef 150g', 'PRD Pe Sea Salt 165g',
       'Smith Crinkle Cut Bolognese 150g', 'Doritos Salsa Mild 300g',
       ], dtype=object)
```

```
In [8]: #Removing special characters
transaction_data_words = transaction_data["PROD_NAME"].str.replace("[^a-z]", '')
```

```
In [9]: pattern = r"([0-9]+(\.g))"
transaction_data_words = transaction_data["PROD_NAME"].str.replace(pattern, '')
```

```
In [10]: transaction_data_words.head()
```

```
Out[10]: 0 Natural Chip Compy SeaSalt
1 CCs Nacho Cheese
2 Smiths Crinkle Cut Chips Chicken
3 Smiths Chip Thiny S/Cream&Onion
4 Kettle Tortilla Chps&ny&lpno Chll
Name: PROD_NAME, dtype: object
```

```
In [11]: words = Counter()
transaction_data_words.str.split().apply(words.update)
print(words)
```

```
Counter({'chips': 49770, 'kettle': 41288, 'i': 35565, 'smiths': 28860, 'salt': 27976, 'cheese': 2789
0, 'pringles': 25102, 'doritos': 24962, 'crinkle': 23960, 'corn': 22065, 'original': 21560, 'cut': 20
754, 'chip': 18645, 'chicken': 18577, 'salsa': 18094, 'sea': 14145, 'thins': 14075, 'sour': 13882, 'c
risps': 12607, 'vinegar': 12402, 'chilli': 12389, 'rip': 11894, 'infusions': 11097, 'supreme': 10963,
'ww': 10320, 'coba': 9850, 'tomato': 9580, 'twisties': 9471, 'twisties': 9464, 'sensa
': 9429, 'old': 9324, 'el': 9324, 'pa': 9324, 'dip': 9324, 'sweet': 7883, 'lime': 7852, 'tomat
o': 7669, 'cream': 7618, 'thiny': 7507, 'tyrellis': 6442, 'and': 6373, 'bbq': 6301, 'tangy': 6332,
'grain': 6272, 'waves': 6272, 'lightly': 6248, 'salted': 6248, 'soy': 6121, 'natural': 6050, 'mild':
6048, 'red': 5885, 'rock': 5885, 'deli': 5885, 'rmai': 4337, 'burger': 4753, 'swt': 4718, 'chives': 4
687, 'honey': 4661, 'nacho': 4658, 'potato': 4647, 'cheezels': 4603, 'ccs': 4551, 'woolworths': 4437,
'mozzarella': 4304, 'basil': 4304, 'pesto': 4304, 'chps&ny&lpno': 3296, 'chilli': 3296, 'swt/chilli':
3242, 'ad&cream': 3269, 'ched': 3268, 'pot': 3257, 'splash': 3252, 'ot': 3252, 'sweetchilli': 3242,
'potato&mx': 3242, 'crinkle': 3233, 'orgni': 3233, 'big': 3233, 'bag': 3233, 'hot': 3229, 'spicy': 312
8, 'cannaberb': 3119, 'fig': 3119, 'barbecue': 3120, 'mexican': 3094, 'jalapeno': 3094, 'lights': 312
8, 'crisps': 3185, 'chip': 3185, 'sweet&spicy': 3177, 'rib': 3174, 'prawn': 3174, 'crackers': 3174, 'so
uthern': 3172, 'strawberry': 3162, 'onion': 3162, 'crn': 3159, 'chives': 3159, 'chps&ny&lpno': 3146,
'chips': 3146, 'chicken': 3146, 'chips': 3146, 'infns': 3144, 'crnchers': 3144, 'gamole': 314
4, 'chps&ny&lpno': 3138, 'sour&cream&herbs': 3134, 'veg': 3134, 'strws': 3134, 'siracha': 3127, 'ch
n': 3125, 'tom': 3125, 'ht': 3125, 'mexicana': 3115, 'seasonedchicken': 3114, 'med': 3114, 'myster
y': 3114, 'flavour': 3114, 'cream&chives': 3105, 'crisps': 3104, 'slit': 3095, 'vingar': 3095, '175g':
3083, 'friedchicken': 3083, 'bbq&maple': 3083, 'rings': 3080, 'chipco': 3010, 'sr': 2984, 'smith': 29
84, 'cheetos': 2977, 'medium': 2879, 'french': 2856, 'snts': 1576, 'whlgrn': 1576, 'ched&matrd': 1
576, 'co': 1572, 'tmat': 1572, 'hbs&spce': 1572, 'vineg': 1550, 'tasty': 1539, 'slow': 1526, 'rat':
1526, 'rock': 1526, 'billy': 1526, 'roast': 1519, 'mac': 1512, 'm': 1512, 'mango': 1507, 'chutny': 15
07, 'papadums': 1507, 'chilli': 1506, 'coconut': 1506, 'snags&sauce': 1503, 'sp': 1498, 'truffle': 14
98, 'chilli&lime': 1495, 'barbecue': 1489, 'stacked': 1487, 'konion&stacked': 1483, 'ch&ny&lpno': 1483,
'chgs': 1479, 'flavour': 1479, 'chips': 1479, 's/cream&onion': 1473, 'pepper': 1473, 'd/style': 1469, 'C
ompy': 1468, 'sea&salt': 1468, 'ornvess': 1468, 'plus': 1468, 'btroot': 1468, 'jam': 1468, 'chili&cr
': 1461, 'hony': 1460, 'chutny': 1460, 'mexzila': 1458, 'steak': 1458, 'chimu&churri': 1458, 'bow':
1454, 'bolognese': 1451, 'puffs': 1448, 'original': 1441, 'salt&d': 1441, 'cut&salt/vineg': 1440, 'onio
n&pb': 1438, 'chiks&garlic': 1434, 'aioli': 1434, 'submits': 1432, 'whlgrn': 1432, 'frch/onin': 143
2, 'ro': 1431, 'ncc': 1419, 'garden': 1419, 'fries': 1418})
```

```
In [12]: words_series = pd.Series(words)
words_series.sort_values(ascending=False)
```

```
Out[12]: chips    49770
kettle    41288
4         35565
Smiths    28860
Salt      27976
...
1432
Pc        1431
NCC       1419
Doritos   1419
Fries     1418
Length: 196, dtype: int64
```

```
In [13]: #There are alot of words with Salsa in them. We want only chips in this dataset so we will remove these
words
print(words_series[words_series.index == "Salsa"])
```

```
Salsa    18094
dtype: int64
```

```
In [14]: pattern = r"([s|a|s|a])"
transaction_data_words = transaction_data.drop(transaction_data[transaction_data["PROD_NAME"].str.contains
(pattern)].index)
transaction_data_words[transaction_data_words["PROD_NAME"].str.contains(pattern)].shape[0]
```

```
Out[14]: 0
```

```
In [15]: #How we will find for outliers and NA values
transaction_data_words.isnull().sum()
```

```
Out[15]: DATE        0
STORE_NBR         0
LYFTY_CARD_NBR    0
TXN_ID            0
PROD_NBR         0
PROD_NAME         0
PROD_QTY          0
TOT_SALES         0
dtype: int64
```

```
In [16]: transaction_data_words.describe()
```

	STORE_NBR	LYFTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_QTY	TOT_SALES
count	24742.000000	2.467420e+05	2.467420e+05	246742.000000	246742.000000	246742.000000
mean	135.051098	1.355310e+05	1.355311e+05	59.351789	1.909062	7.321322
std	76.750708	8.071529e+04	7.814772e+04	33.695428	0.658631	3.077828
min	1.000000	1.000000e+03	1.000000e+03	1.000000	1.000000	1.700000
25%	70.000000	7.001500e+04	6.756925e+04	26.000000	2.000000	5.800000
50%	130.000000	1.303870e+05	1.351638e+05	53.000000	2.000000	7.400000
75%	203.000000	2.030840e+05	2.026538e+05	87.000000	2.000000	8.800000
max	272.000000	2.373711e+06	2.415841e+06	114.000000	200.000000	650.000000

We can clearly see that there are outliers in PROD_QTY and TOT_SALES

```
In [17]: #We can examine the cases where 200 products are bought
transaction_data_words[transaction_data_words["PROD_QTY"] == 200]
```

	DATE	STORE_NBR	LYFTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
69762	2018-08-19	226	226000	226201	4	Dorito Corn Chip Supreme 380g	200	650.0
69763	2019-05-20	226	226000	226210	4	Dorito Corn Chip Supreme 380g	200	650.0

The cases where PROD_QTY is 200, is bought by the same customer since the loyalty card number is the same. Let's see the behavior of this customer, if this customer has other large quantity transactions.

```
In [18]: transaction_data_words[transaction_data_words["LYFTY_CARD_NBR"] == 226000]
```

	DATE	STORE_NBR	LYFTY_CARD_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES
69762	2018-08-19	226	226000	226201	4	Dorito Corn Chip Supreme 380g	200	650.0
69763	2019-05-20	226	226000	226210	4	Dorito Corn Chip Supreme 380g	200	650.0

It looks like this customer has only made two transactions and is not an ordinary retail customer. This customer might be buying for commercial purposes instead. We will remove these two values.

```
In [19]: transaction_data_words = transaction_data_words.drop(transaction_data_words[transaction_data_words["LYFTY_CARD_NBR"]
== 226000].index)
transaction_data_words[transaction_data_words["LYFTY_CARD_NBR"] == 226000].shape[0]
```

```
Out[19]: 0
```

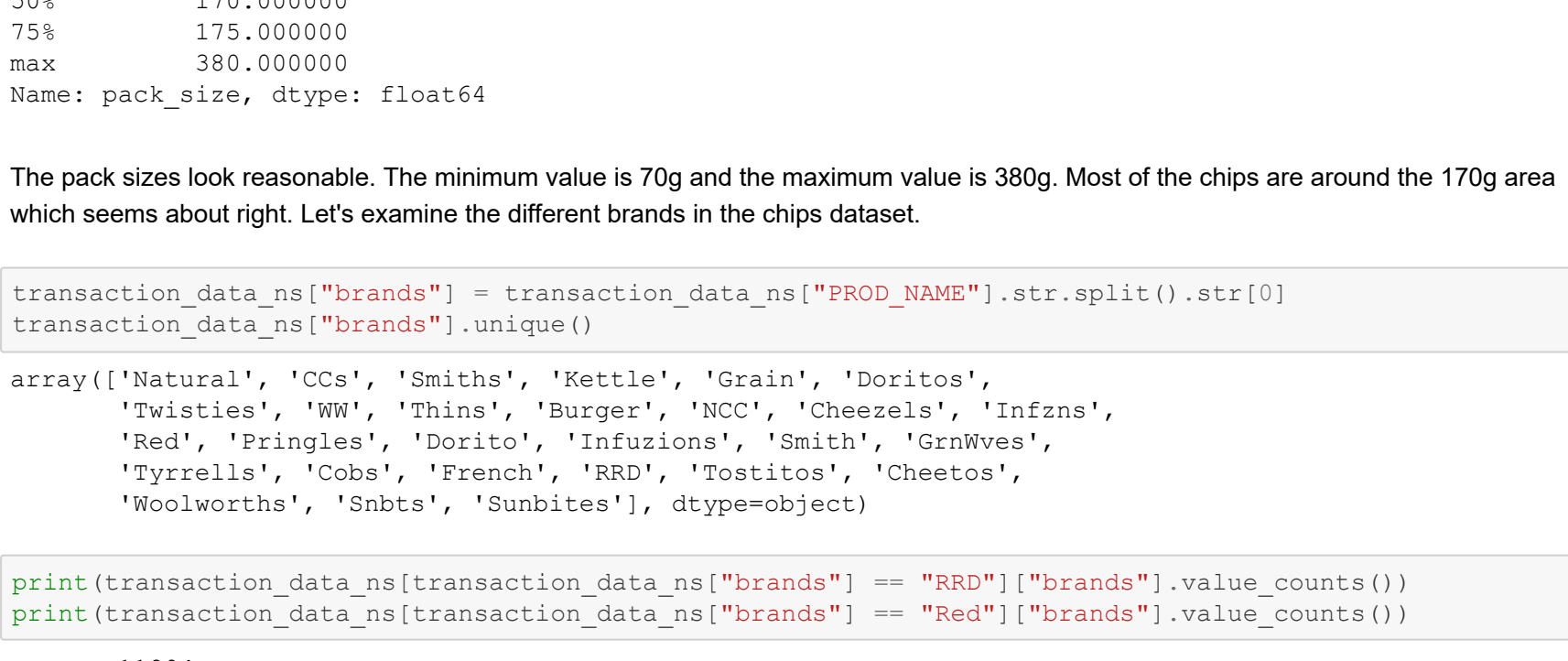
Now we will check the transaction line over time to see if there are any obvious data issues such as missing values

```
In [20]: transaction_per_date = transaction_data_words["DATE"].value_counts()
print(transaction_per_date)

2018-12-24    865
2018-12-23    853
2018-12-22    840
2018-12-19    839
2018-12-20     88
2019-06-24     61
2018-10-18    611
2018-11-25    610
2018-09-22    610
2019-06-13    607
Name: DATE, Length: 364, dtype: int64
```

It seems like there is a missing date since there are only 364 rows

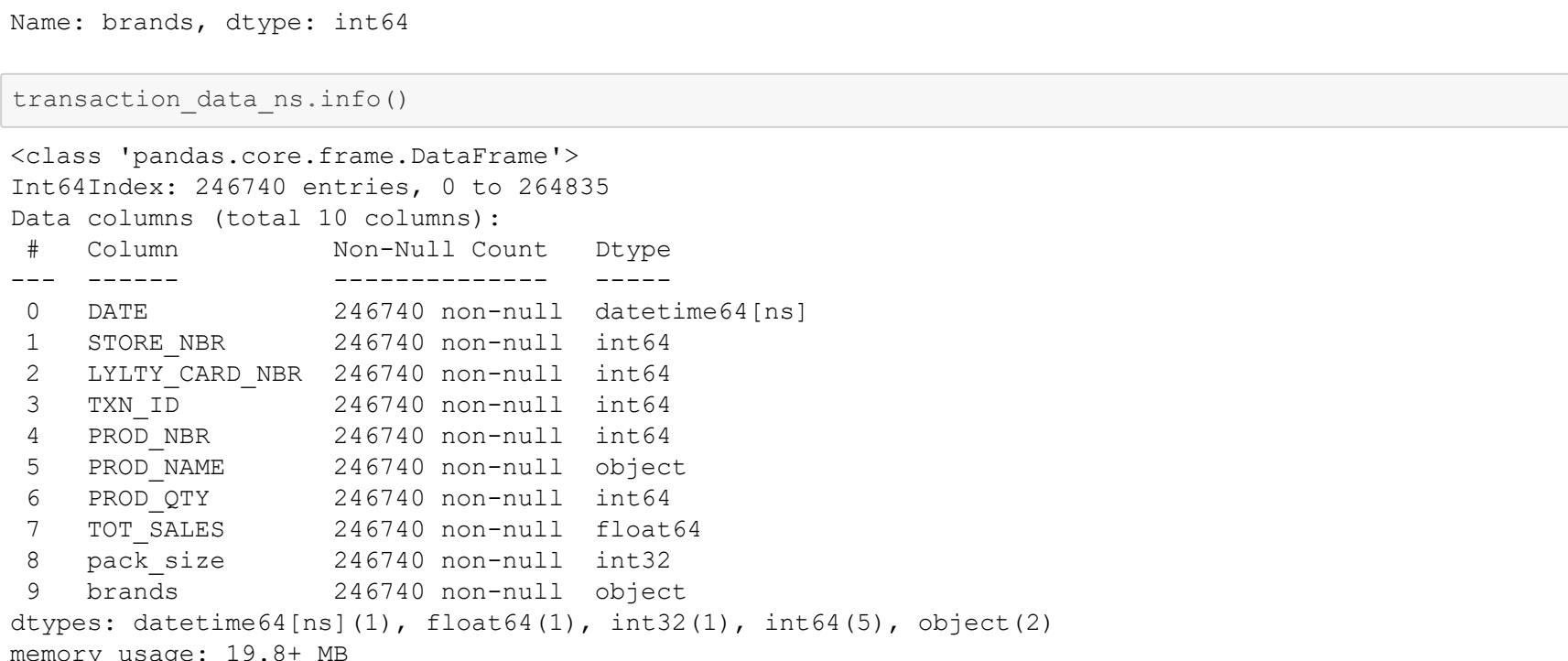
```
In [21]: fig, ax = plt.subplots(figsize=(15,10))
fig, ax = sns.lineplot(transaction_per_date.index, transaction_per_date, ax=ax)
```



It is obvious that there is an outlier between 2018-11 to 2019-01.

```
In [22]: fig1, ax1 = plt.subplots(figsize=(15,10))
fig1, ax1 = sns.lineplot(transaction_per_date.index, transaction_per_date, ax=ax1)
ax1.set_xlim(datetime(2018,12,1), datetime(2019,1,1))
```

```
Out[22]: [(737029.0, 737060.0)]
```



```
In [23]: print(transaction_data_words[transaction_data_words["DATE"] == datetime(2018,12,25)].shape[0])
0
```

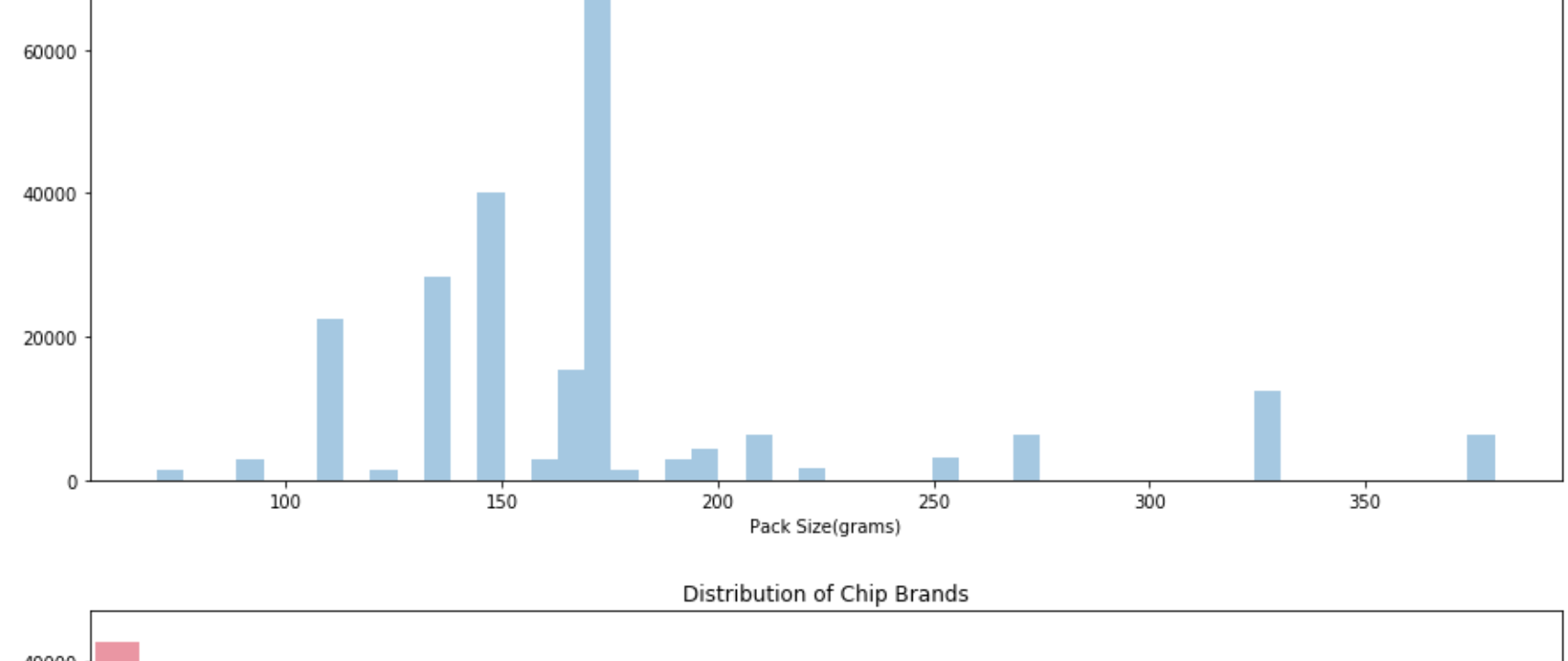
We can see that leading up to christmas, there is a surge in transactions. However, during christmas day, there is 0 purchases. This is due to stores closing during christmas day. This does not count as an outlier. Now we will examine the package size and brands of chips.

```
In [24]: pattern = r"([0-9]+(\.p))"
transaction_data_words["pack_size"] = transaction_data_words["PROD_NAME"].str.extract(pattern).astype(int)
transaction_data_words["pack_size"] = transaction_data_words["pack_size"].unique()
```

```
Out[24]: array([1175, 170, 150, 330, 210, 270, 220, 125, 110, 134, 380, 180, 165,
       135, 250, 200, 160, 190, 90, 70])
```

```
In [25]: fig, ax = plt.subplots(figsize=(15,10))
sns.distplot(transaction_data_words["pack_size"], kde=False)
```

```
Out[25]: <matplotlib.axes._subplots.AxesSubplot at 0x1dc522ba48>
```



```
In [26]: transaction_data_words["pack_size"].describe()
```

count	246740.000000
mean	175.593521
std	59.432118
min	70.000000
25%	150.000000
50%	170.000000
75%	175.000000
max	380.000000
Name: pack_size, dtype: float64	

The pack sizes look reasonable. The minimum value is 70g and the maximum value is 380g. Most of the chips are around the 170g area which seems about right. Let's examine the different brands in the chips dataset.

```
In [27]: transaction_data_words["brands"] = transaction_data_words["PROD_NAME"].str.split().str[0]
```

```
Out[27]: transaction_data_words["brands"].unique()
array(['Natural', 'CCs', 'Smiths', 'Kettle', 'Grain', 'Doritos', 'Twisties', 'WW', 'Thins', 'Burger', 'NCC', 'Cheetos', 'Infns', 'Red', 'Pringles', 'Dorito', 'Infuzions', 'Smith', 'GrnVes', 'Pyrellis', 'Coba', 'French', 'PRD', 'Doritos', 'Chetnos', 'Woolworths', 'Subts', 'Submits'], dtype=object)
```

```
In [28]: print(transaction_data_words[transaction_data_words["brands"] == "RRD"]["brands"].value_counts())
print(transaction_data_words[transaction_data_words["brands"] == "Red"]["brands"].value_counts())
```

RRD	11894
Name: brands, dtype: int64	
Red	4427
Name: brands, dtype: int64	

Some brand names look like they are of the same brands. For example Red and RRD are both Red Rock Deli Chips. We will combine these together.

```
In [29]: transaction_data_words["brands"] = transaction_data_words["brands"].str.lower().str.replace("red", "rrd")
transaction_data_words[transaction_data_words["brands"] == "red"]["brands"].value_counts()
```

```
Out[29]: rrd    16321
Name: brands, dtype: int64
```

```
In [30]: transaction_data_words.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 246740 entries, 0 to 246835
Data columns (total 10 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   DATE        246740 non-null   datetime64[ns]
 1   STORE_NBR   246740 non-null   int64
 2   LIFESTAGE_NBR  246740 non-null   int64
 3   TXN_ID      246740 non-null   int64
 4   PROD_NBR    246740 non-null   int64
 5   PROD_NAME   246740 non-null   object
 6   PROD_QTY    246740 non-null   float64
 7   TOT_SALES   246740 non-null   float64
 8   pack_size   246740 non-null   int32
 9   brands      246740 non-null   object
dtypes: datetime64[ns](1), float64(1), int32(1), int64(5), object(2)
memory usage: 9.8+ MB
```

Data Preparation

Goals:

- Show distributions for brands and pack size
- Combine transaction data and customer_data
- Find for nulls
- Export to CSV

```
In [38]: #Visualizing brand and pack size distributions
brand = pd.Series(transaction_data_words["brands"].value_counts())
fig = plt.figure(figsize=(15,15))
ax1 = fig.add_subplot(2,1,1)
ax2 = fig.add_subplot(2,1,2)
sns.distplot(transaction_data_words["pack_size"], kde=False, ax=ax1, set_title("Distribution of Chip Pack Sizes"))
sns.barplot(x=brand.index, y=brand, ax=ax2).set_title("Distribution of Chip Brands")
ax1.set_xlabel("pack_size", rotation=45)
ax2.set_xlabel("Brands", ylabel="")
ax2.set_ylabel("Count", rotation=45)
```