

# Final Project - Team IMF

In [1]:

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from datetime import date
from datetime import timedelta
import statsmodels.api as sm
from sklearn.linear_model import LinearRegression
import statsmodels.formula.api as smf
from sklearn.metrics import r2_score
from statsmodels.tsa.api import ExponentialSmoothing, SimpleExpSmoothing, Holt
from sklearn.metrics import mean_squared_error
from math import sqrt
from scipy import stats
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix
from sklearn.ensemble import RandomForestClassifier
from sklearn import metrics
from sklearn.model_selection import GridSearchCV
from sklearn import preprocessing
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import precision_recall_curve
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import classification_report
from sklearn.linear_model import LogisticRegression
```

## Summary Stats

In [2]:

```
am=pd.read_csv(r"angels_market.csv")
```

In [3]:

```
am.head()
```

Out[3]:

	vendorID	theme	homeState	carnivals	complaints	est_energy	est_hourly_vol	LL_passh
0	1	Hot Chocolate/Warm Treats	Maine	3	9	57.291961	118	
1	2	Local Artists	Vermont	1	2	39.404898	105	
2	3	Fortune Teller	New Hampshire	5	4	47.175958	94	
3	4	Fried Dough and Pizza	Maine	8	0	58.192568	118	
4	5	craft beer	New Hampshire	7	6	56.657908	102	

In [4]:

```
am.isnull().sum()
```

```
Out[4]: vendorID      0
        theme        0
        homeState    0
        carnivals    0
        complaints   0
        est_energy    0
        est_hourly_vol 0
        LL_passholder 0
        est_hourly_gross 0
        dtype: int64
```

The angels market dataset contains no missing values.

```
In [5]: am.describe()
```

```
Out[5]:
```

	vendorID	carnivals	complaints	est_energy	est_hourly_vol	LL_passholder	est_hourly_gross
<b>count</b>	700.000000	700.000000	700.000000	700.000000	700.000000	700.000000	700.000000
<b>mean</b>	345.375714	5.135714	5.237143	47.501858	110.152857	0.204286	216.54335
<b>std</b>	204.173508	2.202258	4.914525	14.172002	15.903799	0.403467	41.56156
<b>min</b>	1.000000	0.000000	0.000000	3.069903	1.000000	0.000000	4.00000
<b>25%</b>	167.750000	4.000000	0.000000	39.596198	103.000000	0.000000	193.81000
<b>50%</b>	346.500000	5.000000	4.500000	47.955097	110.500000	0.000000	217.49000
<b>75%</b>	521.250000	7.000000	9.000000	57.336190	119.000000	0.000000	242.50250
<b>max</b>	700.000000	13.000000	20.000000	91.567936	147.000000	1.000000	322.57000

```
In [6]: am.groupby('theme')['est_hourly_gross'].describe()
```

```
Out[6]:
```

	count	mean	std	min	25%	50%	75%	max
<b>theme</b>								
<b>3</b>	1.0	6.000000	NaN	6.00	6.0000	6.000	6.0000	6.00
<b>4</b>	1.0	8.000000	NaN	8.00	8.0000	8.000	8.0000	8.00
<b>5</b>	2.0	6.500000	2.121320	5.00	5.7500	6.500	7.2500	8.00
<b>7</b>	2.0	7.500000	0.707107	7.00	7.2500	7.500	7.7500	8.00
<b>8</b>	1.0	7.000000	NaN	7.00	7.0000	7.000	7.0000	7.00
<b>9</b>	1.0	4.000000	NaN	4.00	4.0000	4.000	4.0000	4.00
<b>Canadian Snacks</b>	74.0	221.436892	32.846742	124.75	201.3850	219.380	241.6525	281.79
<b>DIY Ice Sculpture</b>	19.0	222.981053	30.426289	146.03	203.4200	235.810	245.6150	263.95
<b>Fortune Teller</b>	9.0	207.072222	24.840857	175.46	196.1000	200.680	211.0900	263.86
<b>Fried Dough and Pizza</b>	75.0	219.167333	33.469934	144.95	196.2050	217.500	240.0800	291.67
<b>Games Of Chance</b>	85.0	222.085176	32.229983	144.69	203.1500	219.060	242.2700	304.66
<b>Homemade Holiday Gifts</b>	104.0	215.885385	35.286296	81.29	199.5650	217.720	238.9900	286.68
<b>Hot Chocolate/Warm</b>	113.0	214.720354	34.439647	147.03	186.1400	208.810	240.6000	289.43

	count	mean	std	min	25%	50%	75%	max
theme								
<b>Treats</b>								
<b>Local Artists</b>	74.0	224.376216	33.268683	168.56	199.9550	223.995	243.4200	322.57
<b>Local Politician</b>	10.0	222.541000	53.044376	144.15	179.8125	225.930	261.2800	294.42
<b>Maine Tourism Promotion</b>	15.0	215.484000	39.236643	134.64	196.3300	215.290	227.9650	298.84
<b>Specialty Ice Cream</b>	30.0	217.727333	46.960541	137.75	181.8300	212.875	251.2225	305.13
<b>Steaming Hot Cocktails</b>	42.0	218.952381	37.998176	146.69	186.3550	217.465	250.4350	296.16
<b>Video Game/eSports</b>	23.0	217.996522	36.144350	138.03	204.8950	225.590	235.5000	273.78
<b>craft beer</b>	19.0	221.432105	41.563303	146.69	191.8300	228.270	246.9250	285.85

Our first observation was that the dataset contains 6 themes, that are not labelled correctly. This can be due to a data entry error, or vendors with extremely unique themes that did not show up due to scheduling conflicts or insufficient resources to set up their stalls. The latter is more likely the case, because the hourly gross and energy consumption associated with these themes are almost negligible in comparison to the rest of the themes.

Furthermore, we see that Local Artists generate the highest average hourly income among all other vendor themes.

```
In [7]: am['complaints'].value_counts()
```

```
Out[7]: 0      192
        5      46
        3      46
        9      43
        6      43
        7      41
        8      39
        1      39
        4      38
        2      35
       11      26
       10      23
       12      22
       13      18
       14      15
       17      12
       15       9
       16       6
       20       3
       19       2
       18       2
```

Name: complaints, dtype: int64

$192/700=27.42\%$  of the time, customers have had no complaints filed against the vendors in previous carnivals. These would be the safest options for Lobsterland to select for their angels market. However, a glaring limitation with this dataset is the severity of the complaints. If regular complaints can be different from severe complaints, it would be easier for Lobsterland management to select vendors with no serious complaints filed in their past visits.

```
In [8]: am.groupby(by=['carnivals']).mean()
```

```
Out[8]:
```

	vendorID	complaints	est_energy	est_hourly_vol	LL_passholder	est_hourly_gross
<b>carnivals</b>						
0	330.333333	0.000000	53.063056	121.666667	0.333333	262.243333
1	368.310345	6.103448	49.618059	111.137931	0.206897	211.392759
2	387.068182	4.977273	49.382926	112.636364	0.250000	219.491364
3	315.182927	5.451220	48.877076	110.426829	0.134146	223.403293
4	341.622951	5.065574	46.268356	111.581967	0.213115	218.328033
5	359.248120	5.210526	46.254998	108.248120	0.225564	214.606541
6	327.897196	5.859813	48.024971	109.392523	0.196262	213.119720
7	317.192771	5.180723	45.633131	110.373494	0.180723	213.287831
8	379.636364	4.381818	47.858361	110.054545	0.236364	220.763273
9	383.440000	4.360000	51.035565	111.360000	0.200000	223.352000
10	331.166667	4.333333	53.554833	111.333333	0.333333	203.035000
11	367.600000	10.200000	50.731888	106.800000	0.200000	197.976000
12	451.500000	4.000000	54.071423	109.500000	0.000000	175.115000
13	287.750000	5.000000	39.578069	93.500000	0.250000	179.262500

By looking at the different number of carnivals the vendors have appeared in, we can see that the average estimated hourly volume of customers is the highest for a vendor with 0 past carnival apperances, amounting to 121 per hour. Visitors at Lobsterland seem to enjoy exploring new vendors they have not seen at other parks before. The repetitiveness seems to bore customers, as we see how the hourly volume decreases with number of carnival visits.

```
In [9]: am.pivot_table('est_energy', index='carnivals', columns='homeState', margins=True)
```

```
Out[9]:
```

	homeState	2	4	5	6	7	Connecticut	Maine	Massachusetts	New Hampshire	C
<b>carnivals</b>											
0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	67.347419	NaN	NaN	24.
1	NaN	NaN	NaN	NaN	NaN	NaN	70.402830	51.466190	15.123110	51.413808	
2	NaN	NaN	NaN	NaN	NaN	NaN	39.321509	49.078409	43.918733	55.899673	38.
3	NaN	NaN	NaN	NaN	NaN	NaN	49.074666	50.051878	49.420920	43.738546	56.
4	NaN	NaN	NaN	NaN	NaN	NaN	45.274890	45.303136	52.379832	48.519645	54.
5	8.0	9.0	NaN	NaN	5.0	NaN	48.751408	47.293545	45.295418	45.015961	
6	NaN	NaN	8.000000	NaN	NaN	NaN	53.089609	47.941747	37.817052	46.709815	50.
7	NaN	NaN	NaN	NaN	NaN	NaN	46.444979	44.287974	39.836256	47.052452	56.
8	NaN	NaN	4.000000	NaN	5.0	NaN	60.632120	48.504658	44.061156	44.141029	19.

homeState	2	4	5	6	7	Connecticut	Maine	Massachusetts	New Hampshire	C
carnivals										
9	NaN	NaN	NaN	NaN	NaN	NaN	46.064036	NaN	56.671373	74.
10	NaN	NaN	NaN	NaN	NaN	NaN	51.681988	NaN	NaN	
11	NaN	NaN	NaN	NaN	NaN	NaN	54.290327	48.709060	57.043688	
12	NaN	NaN	NaN	NaN	NaN	NaN	50.549629	NaN	57.593216	
13	NaN	NaN	NaN	5.0	NaN	NaN	51.104092	NaN	NaN	
All	8.0	9.0	6.666667	5.0	5.0	49.491742	47.495915	44.172066	47.629391	48.

After creating a pivot table to show the energy consumptions among different homestates, we see that vendors from Vermont have the highest average energy consumption, 51.41% compared to other homestates, whereas vendors from Massachusetts consume the least energy for operating their stalls, at 44.17%. Lobsterland should consider hiring more vendors from MA to save up on utility costs.

## Segmentation and Targeting

In [2]:

```
mf = pd.read_csv(r'maine_families.csv')
```

In [3]:

```
mf
```

Out[3]:

	householdID	total_ppl	own_rent	square_foot	household_income	number_pets	re
0	1	1.0	own	3309	82050.03	1	Aroos
1	2	1.0	own	3814	83077.81	2	Midc
2	3	2.0	rent	2592	91401.41	2	Downeast_Ac
3	4	1.0	own	2628	73048.55	1	Greater Port
4	5	1.0	rent	2442	89145.36	2	Kennebec V
...	...	...	...	...	...	...	
14995	14996	2.0	rent	2802	74859.29	1	Aroos
14996	14997	3.0	own	1906	83083.79	1	Greater Port
14997	14998	2.0	own	3510	109921.74	1	Midc
14998	14999	3.0	rent	2555	47348.86	1	Downeast_Ac
14999	15000	2.0	own	4534	51424.18	0	Greater Port

15000 rows × 10 columns



## Segmentation

Dealing with NaNs with mean values.

In [4]: `mf.isna().sum()`

Out[4]:

householdID	0
total_ppl	75
own_rent	0
square_foot	0
household_income	0
number_pets	0
region	0
entertainment_spend_est	0
travel_spend_est	0
LL_passholder	0
dtype: int64	

In [5]:

```
mean_value=mf['total_ppl'].mean()
mf['total_ppl'].fillna(value=mean_value, inplace=True)
```

In [6]: `mf.isna().sum()`

Out[6]:

householdID	0
total_ppl	0
own_rent	0
square_foot	0
household_income	0
number_pets	0
region	0
entertainment_spend_est	0
travel_spend_est	0
LL_passholder	0
dtype: int64	

**Keep numeric variables only.**

In [7]: `numeric = mf`

In [8]:

```
numeric = numeric.drop('householdID', axis=1)
numeric = numeric.drop('own_rent', axis=1)
numeric = numeric.drop('region', axis=1)
numeric = numeric.drop('LL_passholder', axis=1)
```

In [9]: `numeric`

Out[9]:

	total_ppl	square_foot	household_income	number_pets	entertainment_spend_est	travel_spen
0	1.0	3309	82050.03	1	3189.11	2028.5!
1	1.0	3814	83077.81	2	4175.35	4713.2!
2	2.0	2592	91401.41	2	1814.98	3479.0!
3	1.0	2628	73048.55	1	1945.14	3842.4!
4	1.0	2442	89145.36	2	4410.86	1913.2!
...	...	...	...	...	...	...
14995	2.0	2802	74859.29	1	2878.76	2329.7!
14996	3.0	1906	83083.79	1	2596.40	3456.2!

	total_ppl	square_foot	household_income	number_pets	entertainment_spend_est	travel_spen
<b>14997</b>	2.0	3510	109921.74	1	4836.69	3772.44
<b>14998</b>	3.0	2555	47348.86	1	1148.88	4169.34
<b>14999</b>	2.0	4534	51424.18	0	4458.96	4449.14

15000 rows × 6 columns

## Data scaling

In [10]: `scaling = numeric`

In [11]: `scaling`

Out[11]:

	total_ppl	square_foot	household_income	number_pets	entertainment_spend_est	travel_spen
<b>0</b>	1.0	3309	82050.03	1	3189.11	2028.51
<b>1</b>	1.0	3814	83077.81	2	4175.35	4713.28
<b>2</b>	2.0	2592	91401.41	2	1814.98	3479.01
<b>3</b>	1.0	2628	73048.55	1	1945.14	3842.44
<b>4</b>	1.0	2442	89145.36	2	4410.86	1913.28
...	...	...	...	...	...	...
<b>14995</b>	2.0	2802	74859.29	1	2878.76	2329.78
<b>14996</b>	3.0	1906	83083.79	1	2596.40	3456.21
<b>14997</b>	2.0	3510	109921.74	1	4836.69	3772.44
<b>14998</b>	3.0	2555	47348.86	1	1148.88	4169.34
<b>14999</b>	2.0	4534	51424.18	0	4458.96	4449.14

15000 rows × 6 columns

In [12]:

```

from scipy import stats
scaling['total_ppl'] = stats.zscore(scaling.total_ppl)
scaling['square_foot'] = stats.zscore(scaling.square_foot)
scaling['household_income'] = stats.zscore(scaling.household_income)
scaling['number_pets'] = stats.zscore(scaling.number_pets)
scaling['entertainment_spend_est'] = stats.zscore(scaling.entertainment_spend_est)
scaling['travel_spend_est'] = stats.zscore(scaling.travel_spend_est)

```

In [13]: `scaling`

Out[13]:

	total_ppl	square_foot	household_income	number_pets	entertainment_spend_est	travel_spen
<b>0</b>	-0.714042	0.232960	-0.101357	-0.644446	-0.213953	-1.8
<b>1</b>	-0.714042	0.851909	-0.061449	0.621159	0.588903	1.1

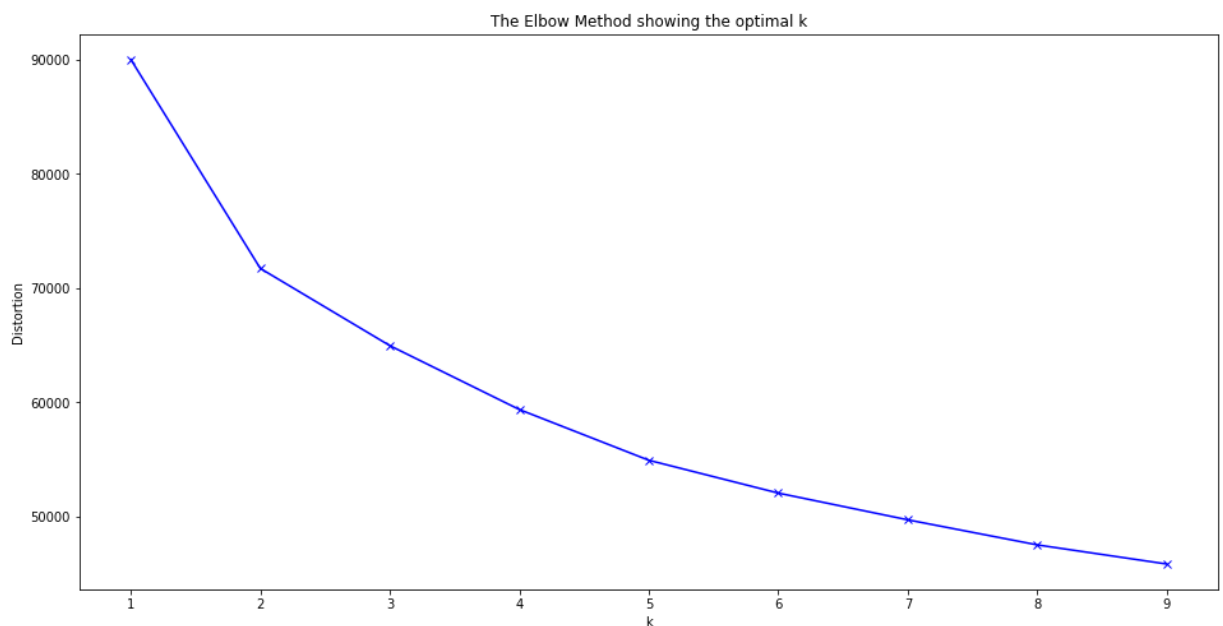
	total_ppl	square_foot	household_income	number_pets	entertainment_spend_est	travel_sper
<b>2</b>	0.515857	-0.645826	0.261750	0.621159	-1.332575	-0.2
<b>3</b>	-0.714042	-0.601703	-0.450877	-0.644446	-1.226617	0.1
<b>4</b>	-0.714042	-0.829672	0.174149	0.621159	0.780622	-1.9
...	...	...	...	...	...	...
<b>14995</b>	0.515857	-0.388441	-0.380568	-0.644446	-0.466596	-1.4
<b>14996</b>	1.745756	-1.486617	-0.061217	-0.644446	-0.696454	-0.2
<b>14997</b>	0.515857	0.479314	0.980880	-0.644446	1.127273	0.0
<b>14998</b>	1.745756	-0.691175	-1.448776	-0.644446	-1.874819	0.5
<b>14999</b>	0.515857	1.734371	-1.290535	-1.910050	0.819778	0.8

15000 rows × 6 columns

In [14]:

```
from sklearn.cluster import KMeans
distortions = []
K = range(1,10)
for k in K:
    kmeanModel = KMeans(n_clusters=k)
    kmeanModel.fit(scoring)
    distortions.append(kmeanModel.inertia_)

plt.figure(figsize=(16,8))
plt.plot(K, distortions, 'bx-')
plt.xlabel('k')
plt.ylabel('Distortion')
plt.title('The Elbow Method showing the optimal k')
plt.show()
```



although the elbow chart shows that 2 might be the best k value for clustering, but I think 4 will be better.

In [15]:

```
kmeanModel = KMeans(n_clusters = 4, random_state = 123)
kmeanModel.fit(scoring)
```



```
scaling['Cluster']=kmeanModel.predict(numeric)
```

```
In [16]: mf['Cluster'] = scaling['Cluster']
```

```
In [17]: numeric['Cluster'].value_counts()
```

```
Out[17]: 0    4107
         2    4001
         3    3715
         1    3177
         Name: Cluster, dtype: int64
```

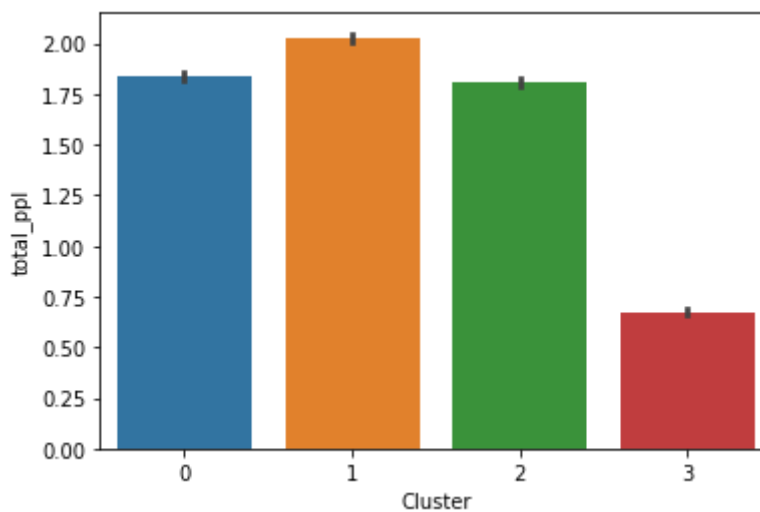
```
In [18]: mf.groupby('Cluster')[['total_ppl',
                                'square_foot',
                                'household_income',
                                'number_pets',
                                'entertainment_spend_est',
                                'travel_spend_est']].mean()
```

```
Out[18]:
```

	total_ppl	square_foot	household_income	number_pets	entertainment_spend_est	travel_spe
<b>Cluster</b>						
<b>0</b>	1.836754	3044.093255	70318.073190	2.230338	2818.574716	3175.2
<b>1</b>	2.026723	3304.224111	108964.347913	1.480327	4726.343009	4398.6
<b>2</b>	1.807080	3019.701325	69843.729708	0.794801	2823.141722	3188.3
<b>3</b>	0.671862	3150.065410	95688.971254	1.506057	3739.469335	4163.4

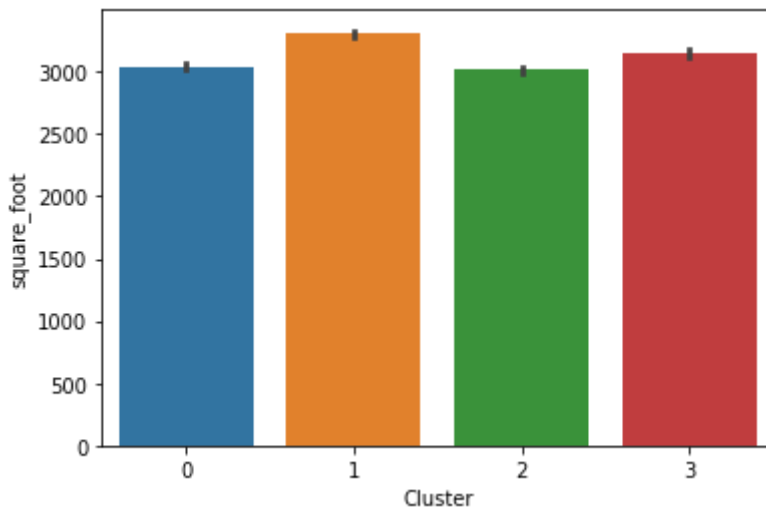
```
In [19]: sns.barplot(x = 'Cluster', y = 'total_ppl', data = mf)
```

```
Out[19]: <AxesSubplot:xlabel='Cluster', ylabel='total_ppl'>
```



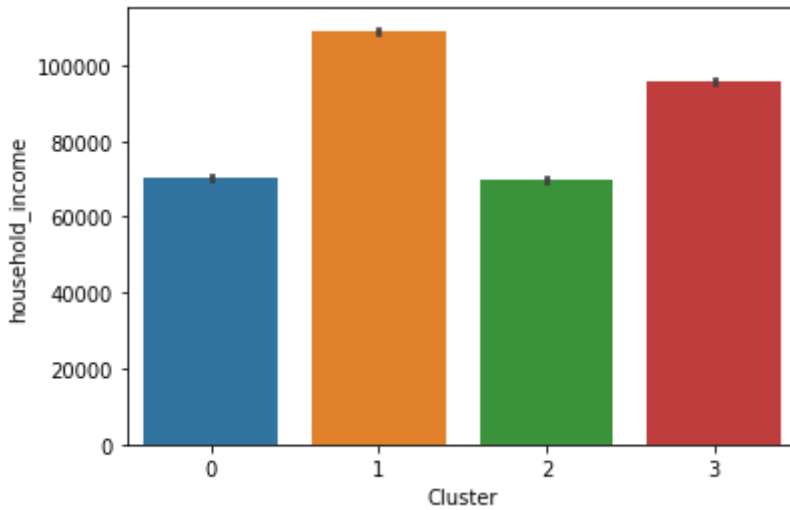
```
In [20]: sns.barplot(x = 'Cluster', y = 'square_foot', data = mf)
```

```
Out[20]: <AxesSubplot:xlabel='Cluster', ylabel='square_foot'>
```



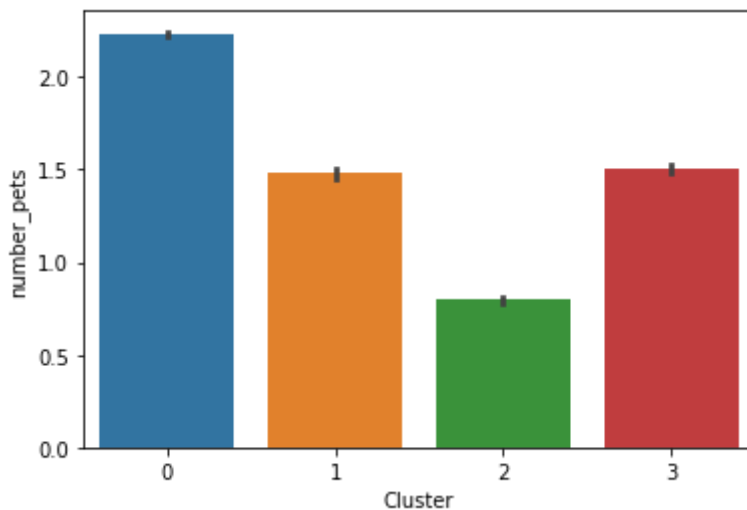
```
In [21]: sns.barplot(x = 'Cluster', y = 'household_income', data = mf)
```

```
Out[21]: <AxesSubplot:xlabel='Cluster', ylabel='household_income'>
```



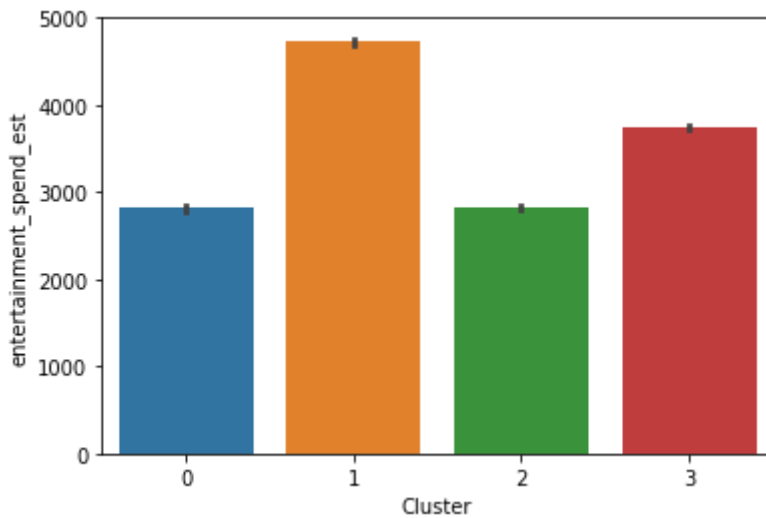
```
In [22]: sns.barplot(x = 'Cluster', y = 'number_pets', data = mf)
```

```
Out[22]: <AxesSubplot:xlabel='Cluster', ylabel='number_pets'>
```



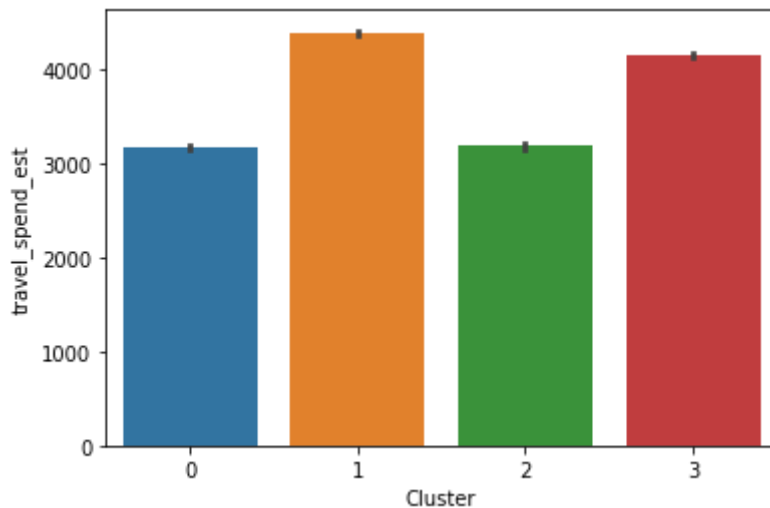
```
In [23]: sns.barplot(x = 'Cluster', y = 'entertainment_spend_est', data = mf)
```

Out[23]: <AxesSubplot:xlabel='Cluster', ylabel='entertainment\_spend\_est'>



In [24]: `sns.barplot(x = 'Cluster', y = 'travel_spend_est', data = mf)`

Out[24]: <AxesSubplot:xlabel='Cluster', ylabel='travel\_spend\_est'>



## Clustering

### Cluster 0: Pet Lover

Cluster 0 is known to me as Pet Lover because they have the largest number of family pets of the four clusters. And they are very willing to invest in their pets, willing to buy bigger homes and spend more on travel for their pets. Sometimes they even sacrifice their entertainment time for their pets.

### Cluster 1: Big Family

I refer to Cluster 1 as a Big Family because they have the highest total household size, total household income and total spending in all areas, which means they are willing to spend, and although they may not have the highest per capita spending, they are undoubtedly the largest customer when it comes to the household as a unit.

### Cluster 2: Low Class

Cluster 2 is what I would call Low Class. The household income and family home size show that this Cluster is the least financially capable of the four clusters. They basically have no pet ties and will keep their spending as low as possible when they go out.

### Cluster 3: Single Elite

I call Cluster 3 Single Elite because they are busy working, some don't even live at home. The size of the home shows their financial strength is strong, and they have the highest per capita income, per capita entertainment spending and per capita travel spending of the four Clusters.

## Targeting

### Cluster 0: Pet Lover

For Cluster 0, since their bond with pets is very deep, I would suggest Lobster Land to introduce some projects where people and pets can play together, or some Pet Special projects such as pet restaurants, pet grooming, etc.

### Cluster 1: Big Family

For Cluster 1, there may be more family members, so Lobster Land can offer them services with a family concept, such as parent-child tickets, or hold entertainment programs that require parents and children to participate together to attract them.

### Cluster 2: Low Class

Even though Cluster 2's spending power is not high, they still have considerable willingness to spend. So for this cluster, what Lobster Land can do is to issue various coupons to them to stimulate them to come and spend.

### Cluster 3: Single Elite

For Cluster 3, their spending power and willingness to spend are very strong. For this cluster, Lobster Land can send them the newly launched high quality entertainment programs to attract them to spend money.

## Conjoint Analysis & Memo Section

```
In [2]: bbq = pd.read_csv("bbq_lake.csv")
```

```
In [3]: bbq.describe()
```

```
Out[3]:
```

	bundleID	avg_rating
count	384.000000	384.000000
mean	192.500000	7.052063
std	110.995495	1.699184
min	1.000000	1.690000
25%	96.750000	5.879500

	bundleID	avg_rating
50%	192.500000	7.121000
75%	288.250000	8.510000
max	384.000000	9.970000

In [4]: `bbq.isnull().sum()`

Out[4]:

bundleID	0
starter	0
maindishI	0
maindishII	0
side	0
dessert	0
avg_rating	0
dtype:	int64

Since there are No NA's in this dataset, we can proceed with dummifying the variables.

In [5]: `bbq.info`

Out[5]:

```
<bound method DataFrame.info of      bundleID      starter      m
aindishI      maindishII  \
0      1  Fried Chicken Tenders      BBQ Brisket      Sausage
1      2  Fried Chicken Tenders      BBQ Brisket      Sausage
2      3  Fried Chicken Tenders      BBQ Brisket      Sausage
3      4  Fried Chicken Tenders      BBQ Brisket      Sausage
4      5  Fried Chicken Tenders      BBQ Brisket      Sausage
..      ...      ...      ...      ...
379    380  Jumbo Shrimp Cocktail  Pork and Brisket Combo  Steak Sampler
380    381  Jumbo Shrimp Cocktail  Pork and Brisket Combo  Steak Sampler
381    382  Jumbo Shrimp Cocktail  Pork and Brisket Combo  Steak Sampler
382    383  Jumbo Shrimp Cocktail  Pork and Brisket Combo  Steak Sampler
383    384  Jumbo Shrimp Cocktail  Pork and Brisket Combo  Steak Sampler

      side      dessert  avg_rating
0      Mac and Cheese      Peach Cobbler      5.81
1      Mac and Cheese  Apple Pie a la Mode      8.93
2      Mashed Potato      Peach Cobbler      6.20
3      Mashed Potato  Apple Pie a la Mode      8.71
4  French Fry Platter      Peach Cobbler      8.24
..      ...      ...      ...
379      Mac and Cheese  Apple Pie a la Mode      5.91
380      Mashed Potato      Peach Cobbler      6.30
381      Mashed Potato  Apple Pie a la Mode      7.65
382  French Fry Platter      Peach Cobbler      8.67
383  French Fry Platter  Apple Pie a la Mode      8.83

[384 rows x 7 columns]>
```

In [6]: `del bbq['bundleID']`

We first imported the dataset and deleted the 'bundleID' variable because it is a unique identifier of the data and has too many levels.

In [7]: `bbq.head()`

Out[7]:

	starter	maindishI	maindishII	side	dessert	avg_rating
0	Fried Chicken Tenders	BBQ Brisket	Sausage	Mac and Cheese	Peach Cobbler	5.81
1	Fried Chicken Tenders	BBQ Brisket	Sausage	Mac and Cheese	Apple Pie a la Mode	8.93
2	Fried Chicken Tenders	BBQ Brisket	Sausage	Mashed Potato	Peach Cobbler	6.20
3	Fried Chicken Tenders	BBQ Brisket	Sausage	Mashed Potato	Apple Pie a la Mode	8.71
4	Fried Chicken Tenders	BBQ Brisket	Sausage	French Fry Platter	Peach Cobbler	8.24

```
In [8]: bbq.columns
```

```
Out[8]: Index(['starter', 'maindishI', 'maindishII', 'side', 'dessert', 'avg_rating'], dtype='object')
```

```
In [9]: bbq01 = pd.get_dummies(bbq, columns=['starter','maindishI','maindishII','side', 'des
```

```
In [10]: bbq01.columns
```

```
Out[10]: Index(['avg_rating', 'starter_Crabcakes and Shrimp',
               'starter_Fried Chicken Tenders', 'starter_Jumbo Shrimp Cocktail',
               'starter_Sticky Chicken Tenders', 'maindishI_BBQ Brisket',
               'maindishI_BBQ Chicken', 'maindishI_Pork and Brisket Combo',
               'maindishI_Pulled Pork', 'maindishII_Beef Short Rib',
               'maindishII_Fajita', 'maindishII_Sausage', 'maindishII_Steak Sampler',
               'side_French Fry Platter', 'side_Mac and Cheese', 'side_Mashed Potato',
               'dessert_Apple Pie a la Mode', 'dessert_Peach Cobbler'],
              dtype='object')
```

I created another instance of the dataframe without drop\_first=True to make sure I have all the levels listed for further analysis.

```
In [11]: bbq1 = pd.get_dummies(bbq, drop_first=True, columns=['starter','maindishI','maindish
```

```
In [12]: bbq1.columns
```

```
Out[12]: Index(['avg_rating', 'starter_Fried Chicken Tenders',
               'starter_Jumbo Shrimp Cocktail', 'starter_Sticky Chicken Tenders',
               'maindishI_BBQ Chicken', 'maindishI_Pork and Brisket Combo',
               'maindishI_Pulled Pork', 'maindishII_Fajita', 'maindishII_Sausage',
               'maindishII_Steak Sampler', 'side_Mac and Cheese', 'side_Mashed Potato',
               'dessert_Peach Cobbler'],
              dtype='object')
```

```
In [13]: X = bbq1[['starter_Fried Chicken Tenders','starter_Jumbo Shrimp Cocktail','starter_S
               'maindishI_Pork and Brisket Combo','maindishI_Pulled Pork','maindishII_Faj
               y = bbq1['avg_rating']
```

```
In [14]: from sklearn.linear_model import LinearRegression
          from sklearn import metrics
```

```
In [15]: regressor = LinearRegression()
          regressor.fit(X, y)
```

Out[15]: `LinearRegression()`

In [16]: `regressor.intercept_`

Out[16]: 6.968645833333337

In [17]: `bbq1.columns`

Out[17]: `Index(['avg_rating', 'starter_Fried Chicken Tenders', 'starter_Jumbo Shrimp Cocktail', 'starter_Sticky Chicken Tenders', 'maindishI_BBQ Chicken', 'maindishI_Pork and Brisket Combo', 'maindishI_Pulled Pork', 'maindishII_Fajita', 'maindishII_Sausage', 'maindishII_Steak Sampler', 'side_Mac and Cheese', 'side_Mashed Potato', 'dessert_Peach Cobbler'], dtype='object')`

In [18]: `coef_df = pd.DataFrame(regressor.coef_, X.columns, columns=['Coefficient'])`  
`coef_df`

Out[18]:

	Coefficient
<b>starter_Fried Chicken Tenders</b>	0.103854
<b>starter_Jumbo Shrimp Cocktail</b>	-0.451771
<b>starter_Sticky Chicken Tenders</b>	0.093333
<b>maindishI_BBQ Chicken</b>	-0.207917
<b>maindishI_Pork and Brisket Combo</b>	1.063646
<b>maindishI_Pulled Pork</b>	0.577187
<b>maindishII_Fajita</b>	-0.498313
<b>maindishII_Sausage</b>	0.342917
<b>maindishII_Steak Sampler</b>	0.232812
<b>side_Mac and Cheese</b>	0.419688
<b>side_Mashed Potato</b>	-0.002500
<b>dessert_Peach Cobbler</b>	-0.739167

After analyzing the coefficients of the linear regression model, built after dummifying the input variables, we can see that the reference levels are as follows:

Starter: Crabcakes and Shrimp

Maindish 1: BBQ Brisket

Maindish 2: Beef Short Rib

Side: French Fry Platter

Dessert: Apple Pie a la Mode

By deep-diving into the strength of the coefficients, we see that for starters, Fried Chicken tenders are the most popular among customers. However, Fried Chicken Tenders cost 3.5

dollars, whereas Sticky chicken tenders costs 2.9 dollars. It is more logical for lobsterland to offer the Sticky chicken tenders as a starter since the vendor cost for sticky chicken is 0.6 dollars lower compared to Fried Chicken tenders, and the difference in coefficients is 11.27 %.  $(0.103854 - 0.093333)/0.093333$ .

In terms of maindish\_1, the 'Pork and Brisket Combo' is the most popular among customers. Another contender for maindish\_1 is 'Pulled Pork', the cost difference between them is only  $6.1 - 6 = 0.1$  dollars. This small difference between their costs is not significant enough to sway our opinion from the obvious coefficient difference of  $(1.063646 - 0.577188)/0.577188$  or 84.28%. So the best choice for maindish\_1 would be 'Pork and Brisket Combo' for 6.1 dollars.

Talking about maindish\_II, 'Sausage' seems to be the most popular choice. 'Steak Sampler' is also a crowd favorite. The cost difference between them is  $5.3 - 4.7 = 0.6$  dollars. The difference in coefficients is  $(0.342917 - 0.232813)/0.232813$  or 47.29% . The cost increase is only  $(5.3 - 4.7)/4.7 = 12.77\%$  . Taking these factors into account, it is more resonable for Lobsterland to offer 'Sausage' as a second maindish despite its higher cost, due to the relative difference in coefficient strength.

For the side, 'Mac and Cheese' is relatively the most popular dish among customers. The french fry platter costs  $0.25 - 0.15 = 0.1$  dollars less compared to mac and cheese, however, based on the relative popularity of mac and cheese, 0.0 compared to 0.419687. Mac and cheese is the best choice here, despite the increased costs. Lobsterland management should select mac and cheese for their preffered side order.

Finally, for desserts, 'The Peach Cobbler' is relatively poor as a dessert choice for customers, compared to 'Apple Pie a la Mode'. The unit cost for this dessert is  $0.9 - 0.6 = 0.3$  dollars higher than Peach Cobbler, however, given the huge negative coefficient associated with Peach Cobbler, Lobsterland should consider having Apple Pie a la Mode as the dessert to make sure its customers are satisfied with the food.

My final recommendation for the menu to Lobsterland would be as follows:

Starter: Sticky Chicken Tenders

Maindish\_1: Pork and Brisket Combo

Maindish\_2: Sausage

Side: Mac and Cheese

Dessert: Apple Pie a la Mode

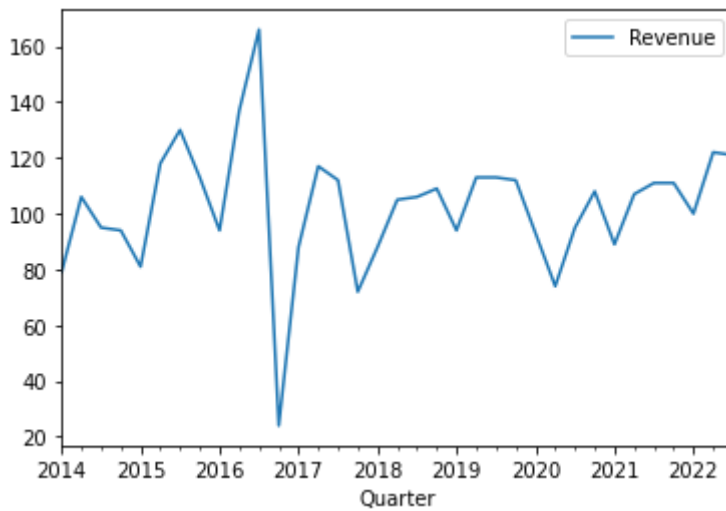
## Forecasting Total Revenue

```
In [2]: df = pd.read_excel('TSQ.xlsx', index_col='Quarter', date_parser = True)
```

```
In [3]: df.plot()
```

```
Out[3]: <AxesSubplot:xlabel='Quarter'>
```



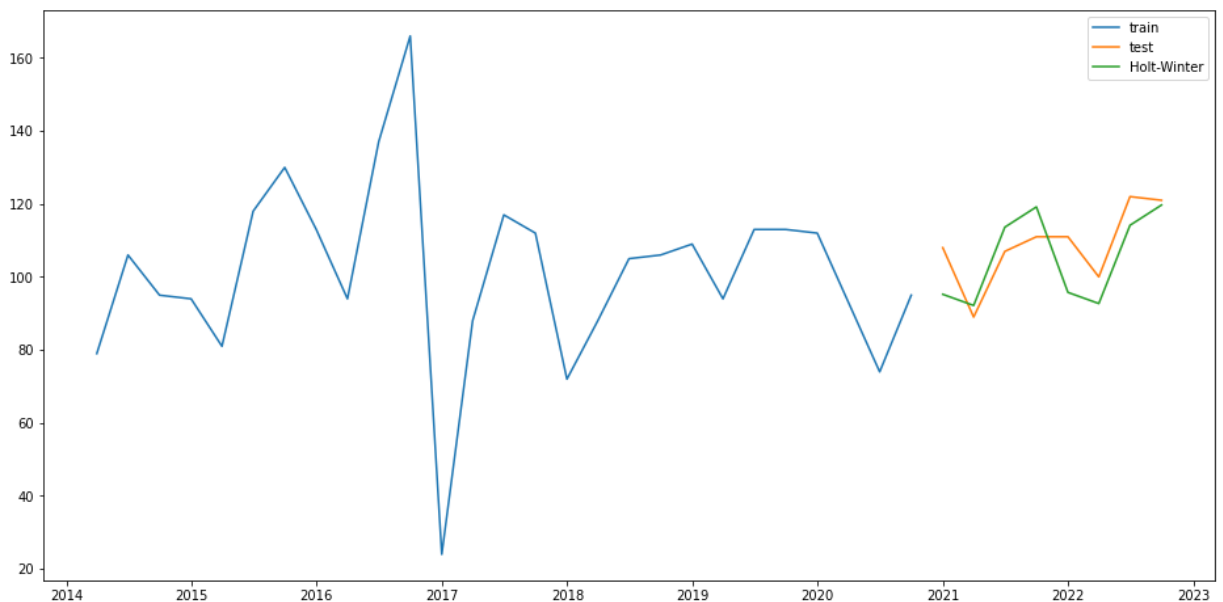


```
In [4]: train = df[0:27]
        test = df[27:]
```

```
In [5]: y_hat_avg = test.copy()
        fit = ExponentialSmoothing(np.asarray(df['Revenue']), seasonal_periods=4, trend='add',
        y_hat_avg['Holt-Winter'] = fit.forecast(len(test))
        Q4_2022 = fit.forecast(1)
        plt.figure(figsize=(16,8))
        plt.plot(train['Revenue'],label='train')
        plt.plot(test['Revenue'],label='test')
        plt.plot(y_hat_avg['Holt-Winter'],label='Holt-Winter')
        plt.legend(loc='best')

        rms = sqrt(mean_squared_error(test['Revenue'],y_hat_avg['Holt-Winter']))
        print(rms)
```

8.892251193700941



```
In [6]: pred_2022 = round(int(y_hat_avg['Holt-Winter'][-3]+ \
        y_hat_avg['Holt-Winter'][-2]+y_hat_avg['Holt-Winter'][-1]+Q4_2
```

```
In [7]: print(f"Total Revenue of TSQ in 2022 is around: {pred_2022} million")
```

Total Revenue of TSQ in 2022 is around: 421 million

In this forecasting case, I tried to use data from Yahoo Finance at first, however, there is only open price, close price and volume. I don't think the revenue can be clarified by those indicators. Therefore, I directly use annual revenue data of TSQ through:

<https://www.macrotrends.net/stocks/charts/TSQ/townsquare-media/revenue>. Based on this data, I also tried Simple Exponential Smoothing and ARIMA, but those methods don't fit well. Compared to annual revenues, quarter revenues have clear pattern and enough data to separate into train and test sets.

## Classification

```
In [2]: cv1 = pd.read_csv(r'carnival_visitors.csv')
cv1.head()
```

```
Out[2]:
```

	householdID	est_inc_USD	est_netw_USD	hhold_field	hhold_oldest	hhold_pax	hhold_youngest
0	23	59245	381931	Govt	48	2	8
1	27	116628	457159	Tech	51	5	21
2	36	65835	394803	Services	50	4	13
3	41	132483	429296	Tech	54	2	11
4	44	83444	488210	Education	51	7	12

## Identify the categorical and numerical variables

```
In [3]: cv1.dtypes
```

```
Out[3]: householdID      int64
est_inc_USD      int64
est_netw_USD      int64
hhold_field      object
hhold_oldest      int64
hhold_pax        int64
hhold_youngest    int64
homeState        object
hhold_car        object
stream_subs      int64
primary          int64
dtype: object
```

So, the above dtypes reflects the total number of categorical and numerical variables in which categorical values are householdID, hhold\_field, homeState, hhold\_car and the numerical variable are est\_inc\_USD, est\_inc\_USD, hhold\_oldest, hhold\_pax, hhold\_youngest, stream\_subs and primary.

```
In [4]: cv1['primary'].value_counts()
```

```
Out[4]: 1    8124
0    6876
Name: primary, dtype: int64
```

So, here we can conclude that the total of 8124 people visited the winter carnival with the purpose of entertainment that can include live performance, concerts, comedy show and team competition and on other hand total 6876 number of people visited the winter carnival with casual consuming mood to enjoy and eat, drink and even engage themselves with shopping and overall, the major purpose of all the visitors is to enjoy, fun and experience cherish moments with their family and loved ones.

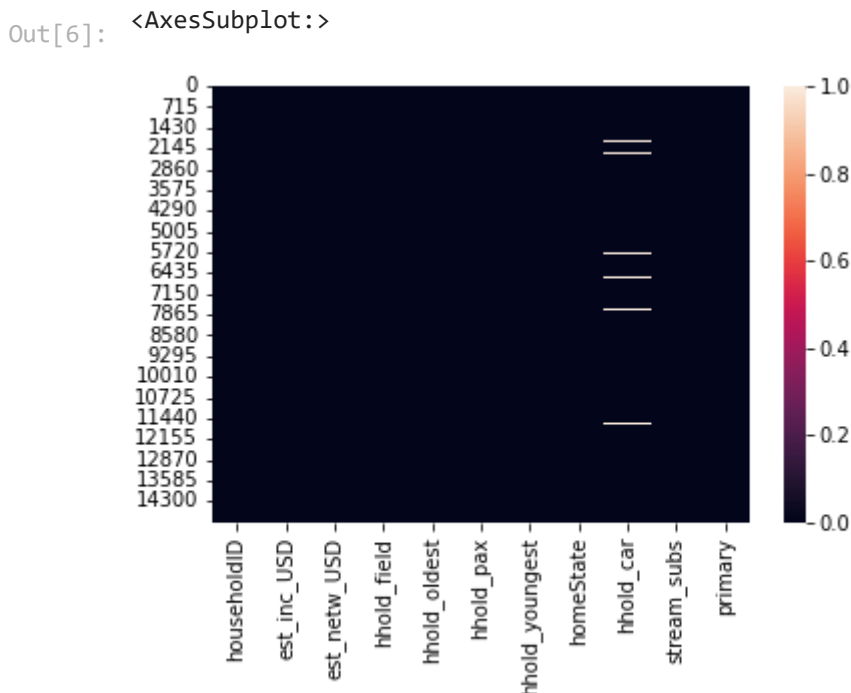
## Dealing with missing values

```
In [5]: cv1.isnull().sum()
```

```
Out[5]: householdID      0
est_inc_USD      0
est_netw_USD      0
hhold_field      0
hhold_oldest      0
hhold_pax      0
hhold_youngest      0
homeState      0
hhold_car      551
stream_subs      0
primary      0
dtype: int64
```

In the above output we can observe that there are total 551 missing values associated with primary source of vehicle transportation and further it is necessary to deal with missing values as the machine learning model will provide an error ahead in the model if we pass NaN values into it and that will disturb the results.

```
In [6]: sns.heatmap(cv1.isnull())
```



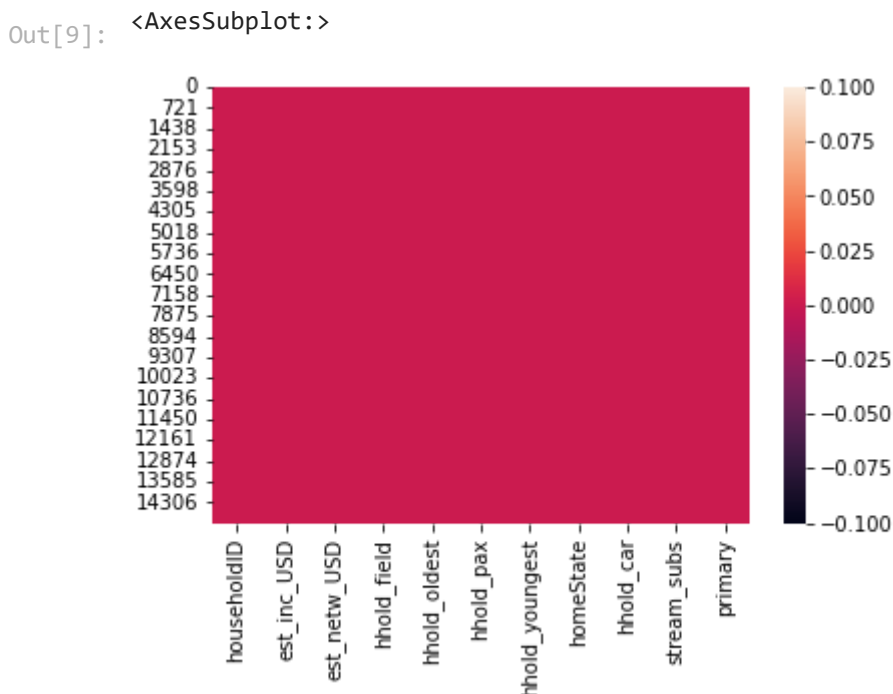
The above heatmap highlights the bunch of lines in the way of hhold\_car that means the visual representation of missing values in the dataset we further try to deal with those missing values with the help of drop isnull function.

```
In [7]: cv1 = cv1.dropna()
```

```
In [8]: cv1.isnull().sum()
```

```
Out[8]: householdID      0
est_inc_USD      0
est_netw_USD     0
hhold_field      0
hhold_oldest     0
hhold_pax        0
hhold_youngest   0
homeState        0
hhold_car        0
stream_subs      0
primary          0
dtype: int64
```

```
In [9]: sns.heatmap(cv1.isnull())
```



This heatmap represent that there are no missing values in the dataset and we are good to go for building the classification model ahead.

## Dummifying

```
In [10]: cv1.columns
```

```
Out[10]: Index(['householdID', 'est_inc_USD', 'est_netw_USD', 'hhold_field',
               'hhold_oldest', 'hhold_pax', 'hhold_youngest', 'homeState', 'hhold_car',
               'stream_subs', 'primary'],
              dtype='object')
```

```
In [11]: cv1 = pd.get_dummies(cv1, drop_first=True, columns=['hhold_field', 'homeState', 'hhold
```

The dummifying the variable is needed to be done as we think that to run the logistic regression model firstly we have dummified the above mentioned categorical variable and on

other hand householdID is okay at its current format as it is unique identifier of every observation in the carnival dataset and numerical data also okay in its current format as we think the only outcome adjusted in logistic regression model is numerical but simultaneously it means that since now one level is dropped so we have to keep this in mind in future.

```
In [12]: print('Shape of dataframe:', cv1.shape)
cv1.head()
```

Shape of dataframe: (14449, 28)

```
Out[12]:
```

	householdID	est_inc_USD	est_netw_USD	hhold_oldest	hhold_pax	hhold_youngest	stream_subs
0	23	59245	381931	48	2	8	2
1	27	116628	457159	51	5	21	3
2	36	65835	394803	50	4	13	3
3	41	132483	429296	54	2	11	3
4	44	83444	488210	51	7	12	3

5 rows × 28 columns

```
In [13]: cv1.columns
```

```
Out[13]: Index(['householdID', 'est_inc_USD', 'est_netw_USD', 'hhold_oldest',
'hhold_pax', 'hhold_youngest', 'stream_subs', 'primary',
'hhold_field_Finance', 'hhold_field_Govt', 'hhold_field_Manufacturing',
'hhold_field_Other', 'hhold_field_Services', 'hhold_field_Tech',
'homeState_Connecticut', 'homeState_Maine', 'homeState_Massachusetts',
'homeState_New Hampshire', 'homeState_New York', 'homeState_Ontario',
'homeState_Quebec', 'homeState_Rhode Island', 'homeState_US_Other',
'homeState_Vermont', 'hhold_car_LuxurySedan', 'hhold_car_Pickup',
'hhold_car_SUV', 'hhold_car_Sedan'],
dtype='object')
```

```
In [14]: X=cv1[['est_inc_USD', 'est_netw_USD', 'hhold_oldest',
'hhold_pax', 'hhold_youngest', 'stream_subs',
'hhold_field_Finance', 'hhold_field_Govt', 'hhold_field_Manufacturing',
'hhold_field_Other', 'hhold_field_Services', 'hhold_field_Tech',
'homeState_Ontario',
'homeState_Quebec', 'homeState_US_Other',
'homeState_Vermont', 'hhold_car_LuxurySedan', 'hhold_car_Pickup',
'hhold_car_SUV', 'hhold_car_Sedan']]
y=cv1['primary']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_stat
```

In order to create data partition we have just taken random part of the data for the purpose of training set and remaining for the test set and we have taken the seed value 74 as the mutual choice of all the group members and we have not taken the householdID while data partition as its again an unique identifier in the carnival dataset which actually don't have strong and valid reason to stay and the most significant chane our team has done to remove the homeState variable as it was providing the disturbing values and the model was not able to highlight the quality output.

```
In [15]: log_reg = sm.Logit(y_train, sm.add_constant(X_train)).fit()
log_reg = sm.Logit(y_test, sm.add_constant(X_test)).fit()
```

Optimization terminated successfully.  
 Current function value: 0.609359  
 Iterations 5  
 Optimization terminated successfully.  
 Current function value: 0.615707  
 Iterations 5

```
In [16]: log_reg.summary()
```

Out[16]:

Logit Regression Results

<b>Dep. Variable:</b>	primary	<b>No. Observations:</b>	5780
<b>Model:</b>	Logit	<b>Df Residuals:</b>	5759
<b>Method:</b>	MLE	<b>Df Model:</b>	20
<b>Date:</b>	Tue, 13 Dec 2022	<b>Pseudo R-squ.:</b>	0.1075
<b>Time:</b>	19:14:34	<b>Log-Likelihood:</b>	-3558.8
<b>converged:</b>	True	<b>LL-Null:</b>	-3987.4
<b>Covariance Type:</b>	nonrobust	<b>LLR p-value:</b>	9.638e-169

	coef	std err	z	P> z	[0.025	0.975]
<b>const</b>	3.8754	0.380	10.199	0.000	3.131	4.620
<b>est_inc_USD</b>	-1.096e-05	1.69e-06	-6.470	0.000	-1.43e-05	-7.64e-06
<b>est_netw_USD</b>	2.569e-06	5.51e-07	4.659	0.000	1.49e-06	3.65e-06
<b>hhold_oldest</b>	-0.1005	0.006	-16.556	0.000	-0.112	-0.089
<b>hhold_pax</b>	0.0393	0.018	2.203	0.028	0.004	0.074
<b>hhold_youngest</b>	-0.0100	0.004	-2.544	0.011	-0.018	-0.002
<b>stream_subs</b>	0.2397	0.017	14.123	0.000	0.206	0.273
<b>hhold_field_Finance</b>	-0.7556	0.111	-6.782	0.000	-0.974	-0.537
<b>hhold_field_Govt</b>	-0.9747	0.097	-10.053	0.000	-1.165	-0.785
<b>hhold_field_Manufacturing</b>	-0.5779	0.177	-3.259	0.001	-0.925	-0.230
<b>hhold_field_Other</b>	-0.3646	0.200	-1.819	0.069	-0.757	0.028
<b>hhold_field_Services</b>	0.2330	0.112	2.080	0.038	0.013	0.453
<b>hhold_field_Tech</b>	-0.4378	0.114	-3.833	0.000	-0.662	-0.214
<b>homeState_Ontario</b>	-0.0615	0.123	-0.501	0.616	-0.302	0.179
<b>homeState_Quebec</b>	-0.0130	0.119	-0.110	0.912	-0.245	0.219
<b>homeState_US_Other</b>	-0.1375	0.141	-0.973	0.331	-0.415	0.139
<b>homeState_Vermont</b>	-0.0160	0.096	-0.167	0.868	-0.204	0.172
<b>hhold_car_LuxurySedan</b>	0.3469	0.115	3.029	0.002	0.122	0.571
<b>hhold_car_Pickup</b>	0.0126	0.133	0.095	0.924	-0.247	0.273
<b>hhold_car_SUV</b>	0.1783	0.108	1.647	0.100	-0.034	0.390

<b>hhold_car_Sedan</b>	0.3935	0.118	3.345	0.001	0.163	0.624
------------------------	--------	-------	-------	-------	-------	-------

After building the model the above summary indicates some high p-value of some the variable in which the homestate\_quebec have maximum p-value which can brings lots of reson behind this but the most impactable is the different country as it's hard to travel all the way from canada or the time, fuel, resorces would be extra as compared to people live nearby and finally it does not give much statistical significance to the model.

## Shape and look of training and testing sets

```
In [17]: print('Shape of training feature:', X_train.shape)
print('Shape of testing feature:', X_test.shape)
print('Shape of training label:', y_train.shape)
print('Shape of training label:', y_test.shape)
```

```
Shape of training feature: (8669, 20)
Shape of testing feature: (5780, 20)
Shape of training label: (8669,)
Shape of training label: (5780,)
```

The shape defines the clear shape of both training as well as testing set that ultimately shows the twenty column in both the set but with different numbers of columns.

## Predicting the model

```
In [18]: classifier = KNeighborsClassifier(n_neighbors=5)
classifier.fit(X_train, y_train)
```

```
Out[18]: KNeighborsClassifier()
```

```
In [19]: knn = KNeighborsClassifier()
knn.fit(X_train, y_train)
predictions = knn.predict(X_test)
print(accuracy_score(y_test, predictions))
print(confusion_matrix(y_test, predictions))
print(classification_report(y_test, predictions))
```

```
0.5377162629757786
```

```
[[1181 1475]
 [1197 1927]]
```

	precision	recall	f1-score	support
0	0.50	0.44	0.47	2656
1	0.57	0.62	0.59	3124
accuracy			0.54	5780
macro avg	0.53	0.53	0.53	5780
weighted avg	0.53	0.54	0.53	5780

In this model we can identify that the accuracy score is 54% and sensitivity rate is 62% with the specificity rate is 44% through which we can conclude that the balanced accuracy of the model is sensitivity+specificity/2 which would be 53%.

## Predicting through decision tree classification

```
In [20]: dtc = DecisionTreeClassifier()
dtc.fit(X_train, y_train)
predictions = dtc.predict(X_test)
print(accuracy_score(y_test, predictions))
print(confusion_matrix(y_test, predictions))
print(classification_report(y_test, predictions))
```

```
0.6200692041522491
```

```
[[1568 1088]
```

```
 [1108 2016]]
```

	precision	recall	f1-score	support
0	0.59	0.59	0.59	2656
1	0.65	0.65	0.65	3124
accuracy			0.62	5780
macro avg	0.62	0.62	0.62	5780
weighted avg	0.62	0.62	0.62	5780

We have predicted the model again with decision tree classifier as it works for both continuous as well as categorical output variables and also this suggest our team to determine the potential and primary purpose of the visitors to visit in winter carnival for either consume or entertain purpose and in this model we can identify that the accuracy score is 61% and sensitivity rate is 63% with the specificity rate is 58% through which we can conclude that the balanced accuracy of the model is sensitivity+specificity/2 which would be 60.5%.

## Accuracy of train and test set

```
In [21]: y_predict = classifier.predict(X_test)
logmodel = LogisticRegression(max_iter=500)
logmodel.fit(X_train, y_train)
```

```
Out[21]: LogisticRegression(max_iter=500)
```

```
In [22]: predictions = logmodel.predict(X_test)
accuracy_score(y_test, predictions)
```

```
Out[22]: 0.5826989619377163
```

```
In [23]: predictions = logmodel.predict(X_train)
print(metrics.accuracy_score(y_train, predictions))
```

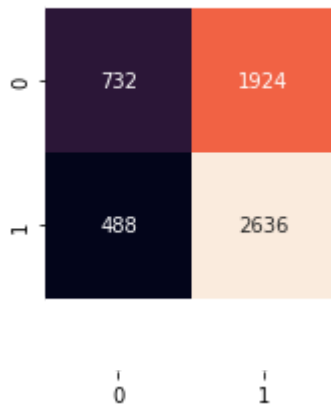
```
0.5734225400853616
```

In this section we can observe that the accuracy score of testing set is 0.5826 and on other hand the accracy score of training set is 0.5734.

```
In [24]: %matplotlib inline
predictions = logmodel.predict(X_test)
mat = confusion_matrix(y_test, predictions)
sns.heatmap(mat, fmt='g', square=True, annot=True, cbar=False)
a, b = plt.ylim()
a += 0.5
b -= 0.5
```



```
plt.ylim(a, b)
```



So, what exactly does this figure tell us about the performance of our model? Looking along the diagonal of the confusion matrix, let's pay attention to the numbers 732 and 2636. The number 732 corresponds to the number of visitors who were correctly predicted by the model who visit the winter carnival with mood of consuming the food and drinks and enjoy the evening, meaning they will contribute more towards carnival merchandise and snack stores that can generate stronger revenue model and the number 2636 corresponds to the number of visitors that the model correctly predicted to visit the winter carnival with entertain purpose and experience the amazing live performance, comedy shows and competitions that will create thrill and excitement among the visitors and probably they'll refer and recommend their friends and neighbours to visit the winter carnival that will also simultaneously increase the visitors traffic and profit margin of the park. The and ultimately with the light of confusion matrix model and the p-value of all the variables suggest that visitors have more desire to visit winter carnival with entertainment purpose and to experience the memorable and cherish moments while watching the thriller competitions and live performances in winter carnival and apart from this winter carnival could improve the food section of the park by introducing the bunch of discount offers and bundle offers that attracts more customer and convince them to spend more money in the park but there can be actually several more reasons why people visit the winter carnival as they live nearby and getting bored sticking in front of the screen and think to visit park again or they might have some promotional discount offer or one day pass and many more but in nutshell winter carnival have different customer segmentations and priority is to increase the popularity and revenue model of the park.

## A/B Testing

```
In [2]: sp = pd.read_csv(r'snowmobile_pics.csv')
        sp
```

```
Out[2]:
```

	recipient	pic_seen	site_duration	spend	register
0	1	Racers in Action	18.20	16.60	0
1	2	Starting Line	28.61	15.30	0
2	3	Sharp Turn	10.90	16.32	1
3	4	Sharp Turn	11.30	22.62	0

	recipient	pic_seen	site_duration	spend	register
4	5	Racers in Action	19.70	17.30	0
...	...	...	...	...	...
3395	3396	Sharp Turn	11.80	19.12	1
3396	3397	Starting Line	23.61	16.80	0
3397	3398	Sharp Turn	10.90	18.82	1
3398	3399	Starting Line	14.71	20.50	0
3399	3400	Racers in Action	22.90	17.50	0

3400 rows × 5 columns

In [3]:

sp.groupby('pic\_seen').describe()

Out[3]:

	recipient							site_duration		..	
	count	mean	std	min	25%	50%	75%	max	count	mean	..
pic_seen											
Racers in Action	1110.0	1693.945946	993.119415	1.0	820.50	1688.5	2573.50	3400.0	1110.0	22.949189	..
Sharp Turn	1142.0	1734.473730	977.916270	3.0	870.25	1785.0	2560.50	3398.0	1142.0	10.975394	..
Starting Line	1148.0	1673.040941	974.007851	2.0	847.75	1644.5	2511.75	3399.0	1148.0	24.294059	..

3 rows × 32 columns



In [4]:

sp.groupby('pic\_seen').mean()

Out[4]:

	recipient	site_duration	spend	register
pic_seen				
Racers in Action	1693.945946	22.949189	16.781892	0.440541
Sharp Turn	1734.473730	10.975394	18.606778	0.348511
Starting Line	1673.040941	24.294059	14.016289	0.341463

Set alpha threshold to be .05 for all comparisons.

Based on the stats shown above, I claim that each type of pic is best for a specific kpi.

## Situation I

lobster Land prioritizing the kpi "register"

i. Pic: Racers in Action vs. Pic: Sharp Turn

The null hypothesis  $H_0$  is that Pic: Racers in Action and Pic: Sharp Turn are equally effective at "register".

```
In [5]: t, p = stats.ttest_ind(sp.loc[sp['pic_seen'] == 'Racers in Action', 'register'].values, sp.loc[sp['pic_seen'] == 'Sharp Turn', 'register'].values)
print('t-value is equal to '+str(t)+',', ' p-value is equal to ' + str(p))
```

t-value is equal to 4.483983175631638, p-value is equal to 7.69613709155986e-06

Since the p-value for our t-test is less than .05, we reject the null hypothesis that Pic: Racers in Action and Pic: Sharp Turn are equally effective at "register".

Furthermore, 44.05% recipients who received "Racers in Action" registered but only 34.85% recipients who received "Sharp Turn" registered.

Therefore, we can conclude that "Racers in Action" is more effective than "Sharp Turn" at the kpi "register".

ii. Pic: Racers in Action vs. Pic: Starting Line

The null hypothesis  $H_0$  is that Pic: Racers in Action and Pic: Starting Line are equally effective at "register".

```
In [6]: t, p = stats.ttest_ind(sp.loc[sp['pic_seen'] == 'Racers in Action', 'register'].values, sp.loc[sp['pic_seen'] == 'Starting Line', 'register'].values)
print('t-value is equal to '+str(t)+',', ' p-value is equal to ' + str(p))
```

t-value is equal to 4.844354466304721, p-value is equal to 1.357059422760036e-06

Since the p-value for our t-test is less than .05, we reject the null hypothesis that Pic: Racers in Action and Pic: Starting are equally effective at "register".

Furthermore, 44.05% recipients who received "Racers in Action" registered but only 34.15% recipients who received "Starting Line" registered.

Therefore, we can conclude that "Racers in Action" is more effective than "Sharp Turn" at the kpi "register".

iii. Pic: Sharp Turn vs. Pic: Starting Line

The null hypothesis  $H_0$  is that Pic: Sharp Turn and Pic: Starting Line are equally effective at "register".

```
In [7]: t, p = stats.ttest_ind(sp.loc[sp['pic_seen'] == 'Sharp Turn', 'register'].values, sp.loc[sp['pic_seen'] == 'Starting Line', 'register'].values)
print('t-value is equal to '+str(t)+',', ' p-value is equal to ' + str(p))
```

t-value is equal to 0.3546036791692854, p-value is equal to 0.7229192023249931

Since the p-value for our t-test is more than .05, we fail to reject the null hypothesis that Pic: Sharp Turn and Pic: Starting are equally effective at "register".

Furthermore, 34.85% recipients who received "Racers in Action" registered and 34.15% recipients who received "Starting Line" registered. These two percentages are indeed pretty close.

Therefore, we can conclude that Pic: Sharp Turn and Pic: Starting Line are equally effective at the kpi "register".

In summary, for the kpi "register", "Racers in Action" > "Sharp Turn" = "Starting Line"

## Situation II

Lobster Land prioritizing the kpi site\_duration

i. Pic: Racers in Action vs. Pic: Sharp Turn

The null hypothesis  $H_0$  is that Pic: Racers in Action and Pic: Sharp Turn are equally effective at "site\_duration".

```
In [8]: t, p = stats.ttest_ind(sp.loc[sp['pic_seen'] == 'Racers in Action', 'site_duration'],
                             print('t-value is equal to '+str(t)+',', ' p-value is equal to ' + str(p))
```

t-value is equal to 178.43526393191596, p-value is equal to 0.0

Since the p-value for our t-test is less than .05, we reject the null hypothesis that Pic: Racers in Action and Pic: Sharp Turn are equally effective at "site\_duration".

Also, since the mean of site\_duration for Racers in Action is 22.95 and the mean of site\_duration for Sharp Turn is 10.98, we can conclude that Racers in Action is more effective than Sharp Turn at site\_duration.

ii. Pic: Racers in Action vs. Pic: Starting Line

The null hypothesis  $H_0$  is that Pic: Racers in Action and Pic: Starting Line are equally effective at "site\_duration".

```
In [9]: t, p = stats.ttest_ind(sp.loc[sp['pic_seen'] == 'Racers in Action', 'site_duration'],
                             print('t-value is equal to '+str(t)+',', ' p-value is equal to ' + str(p))
```

t-value is equal to -10.001841948108055, p-value is equal to 5.917256397701659e-23

Since the p-value for our t-test is less than .05, we reject the null hypothesis that Pic: Racers in Action and Pic: Starting Line are equally effective at "site\_duration".

Also, since the mean of site\_duration for Racers in Action is 22.95 and the mean of site\_duration for Starting Line is 24.29, we can conclude that Racers in Action is less effective than Starting Line at site\_duration.

iii. Pic: Sharp Turn vs. Pic: Starting Line

The null hypothesis  $H_0$  is that Pic: Sharp Turn and Pic: Starting Line are equally effective at "site\_duration".

```
In [10]: t, p = stats.ttest_ind(sp.loc[sp['pic_seen'] == 'Sharp Turn', 'site_duration'],
                                print('t-value is equal to '+str(t)+',', ' p-value is equal to ' + str(p))
```

t-value is equal to -112.41074656435279, p-value is equal to 0.0

Since the p-value for our t-test is less than .05, we reject the null hypothesis that Pic: Sharp Turn and Pic: Starting Line are equally effective at "site\_duration".

Also, since the mean of site\_duration for Sharp Turn is 10.98 and the mean of site\_duration for Starting Line is 24.29, we can conclude that Sharp Turn is less effective than Starting Line at site\_duration.

In summary, for the kpi site\_duration, Starting Line > Racers in Action > Sharp Turn

## Situation III

Lobster Land prioritizing the kpi "spend"

i. Pic: Racers in Action vs. Pic: Sharp Turn

The null hypothesis  $H_0$  is that Pic: Racers in Action and Pic: Sharp Turn are equally effective at "spend".

```
In [11]: t, p = stats.ttest_ind(sp.loc[sp['pic_seen'] == 'Racers in Action', 'spend'].values,
print('t-value is equal to '+str(t)+',', ' p-value is equal to ' + str(p))
```

t-value is equal to -23.859305670179317, p-value is equal to 3.8873183005703336e-106

Since the p-value for our t-test is less than .05, we reject the null hypothesis that Pic: Racers in Action and Pic: Sharp Turn are equally effective at "spend".

Also, since the mean of spend for Racers in Action is 16.78 and the mean of spend for Sharp Turn is 18.61, we can conclude that Racers in Action is less effective than Sharp Turn at spend.

ii. Pic: Racers in Action vs. Pic: Starting Line

The null hypothesis  $H_0$  is that Pic: Racers in Action and Pic: Starting Line are equally effective at "spend".

```
In [12]: t, p = stats.ttest_ind(sp.loc[sp['pic_seen'] == 'Racers in Action', 'spend'].values,
print('t-value is equal to '+str(t)+',', ' p-value is equal to ' + str(p))
```

t-value is equal to 41.02687256136666, p-value is equal to 1.7008147207546044e-248

Since the p-value for our t-test is less than .05, we reject the null hypothesis that Pic: Racers in Action and Pic: Starting Line are equally effective at "spend".

Also, since the mean of spend for Racers in Action is 16.78 and the mean of spend for Starting Line is 14.02, we can conclude that Racers in Action is more effective than Starting Line at spend.

In summary, for the kpi "spend", we have Sharp Turn > Racers in Action > Starting Line

Putting every kpi comparison together, we have:

For the kpi "register", "Racers in Action" > "Sharp Turn" = "Starting Line"

For the kpi "site\_duration", Starting Line > Racers in Action > Sharp Turn

For the kpi "spend", we have Sharp Turn > Racers in Action > Starting Line

## Recommendation to Lobster Land

If Lobster Land can prioritize a specific kpi, it be be easier for it to choose what picture to use based on above summary.

However, if Lobster Land has to choose a pic to use without targetting a specific kpi, I'd recommendation it to choose "Racers in Action" as "Racers in Action" performs the best in general.

## Conclusions

In this project, our group first subdivided the visitors into four clusters and developed corresponding targeting strategies for each type of visitors. Then, using conjoint analysis, the most popular menu items in the Lobster Land Barbeque Tent were identified. Then we made a forecast based on historical data and concluded that Town Square Media's total revenue at the end of 2022 would be about \$421 million. We then built a classification model to predict whether visitors would prefer to spend on snacks or shows, and in light of the success of the Qingdao Beer Festival in China, we tried to learn from this example to replicate the success at Lobster Land. Finally, we conducted a series of A/B tests and concluded that "Racers in Action" was the most popular of the three photos.

We believe that all of the above analyses are very useful for Lobster Land, as they allow us to better understand our customers and adjust our business strategies to maximize revenue. For example, identifying customer clusters and implementing specific marketing approaches for each cluster will undoubtedly lead to more efficient conversion rates and improved customer retention. Identifying Barbeque Tent's most popular menu items will help the restaurant achieve optimal costs and reduce waste when planning purchases.