

# Starbucks Data Analysis Project

The problem that I chose to solve was to build a model that predicts whether a customer will respond to a Starbucks offer. The data set that I will be using for my analysis is simulated data that mimics customer behavior on the Starbucks rewards mobile app. Once every few days, Starbucks sends out an offer to users of the mobile app. An offer can be merely an advertisement for a drink or an actual offer such as a discount or BOGO (buy one get one free). Some users might not receive any offer during certain weeks. Not all users receive the same offer, and this is the challenge to solve with this data set. My task is to combine transaction, demographic and offer data to determine which demographic groups respond best to which offer type. This data set is a simplified version of the real Starbucks app because the underlying simulator only has one product whereas Starbucks actually sells dozens of products. Every offer has a validity period before the offer expires. As an example, a BOGO offer might be valid for only 5 days. In the data set informational offers have a validity period even though these ads are merely providing information about a product for purchase, for example, if an informational offer has 7 days of validity, you can assume the customer is feeling the influence of the offer for 7 days after receiving the advertisement. I have extracted transactional data showing user purchases made on the app including the timestamp of purchase and the amount of money spent on a purchase. This transactional data also has a record for each offer that a user receives as well as a record for when a user actually views the offer. There are also records for when a user completes an offer. I also kept in mind that someone using the app might make a purchase through the app without having received an offer or seen an offer. To give an example, a user could receive a discount offer buy 10 dollars get 2 off on Monday. The offer is valid for 10 days from receipt. If the customer accumulates at least 10 dollars in purchases during the validity period, the customer completes the offer. However, there are a few things to watch out for in this data set. Customers do not opt into the offers that they receive; in other words, a user can receive an offer, never actually view the offer, and still complete the offer. For example, a user might receive the "buy 10 dollars get 2 dollars off" offer, but the user never opens the offer during the 10 day validity period. The customer spends 15 dollars during those ten days. There will be an offer completion record in the data set; however, the customer was not influenced by the offer because the customer never viewed the offer. Aim is to build a machine learning model that could help predicting which customer will actually use the offer.

The data is contained in three files:

portfolio.json - containing offer ids and meta data about each offer (duration, type, etc.), profile.json - demographic data for each customer. transcript.json - records for transactions, offers received, offers viewed, and offers completed. Here is the schema and explanation of each variable in the files: # portfolio.json

id (string)- offer id, offer type (string)- the type of offer is BOGO, discount, informational, difficulty (int)-the minimum required to spend to complete an offer, reward (int)-the reward is given for completing an offer, duration (int)-time for the offer to be open, in days, channels (list of strings). # profile.json

age (int)-age of the customer, became\_member\_on (datetime)-the date when customer created an app account, gender (str)-gender of the customer (note some entries contain "Q" for other than M or F), id (str)-customer ID, income (float)-customer's income, # transcript.json event (obj)-record describing (is transaction, offer received, offer viewed, etc.), person (obj)-customer id, time (int) in hours since the start of the data. The data begins at time=0, value=(dict of strings)-another offer id or transaction amount depending on the record.

Download dataset <https://www.kaggle.com/blacklitter/starbucks-app-customer-reward-program-data>

```
In [9]: import pandas as pd
import numpy as np
import math
import json
import matplotlib.pyplot as plt
import seaborn as sns

# read in the json files
portfolio = pd.read_json('portfolio.json', orient='records', lines=True)
profile = pd.read_json('profile.json', orient='records', lines=True)
transcript = pd.read_json('transcript.json', orient='records', lines=True)
```

```
In [5]: portfolio.head()

Out[5]:
```

reward	difficulty	channels	duration	offer_type	id
0	10	[email, mobile, social]	10	5	bogo
1	10	[web, email, mobile, social]	5	7	bogo
2	0	[web, email, mobile]	4	informational	3f207d87b143eaa3ce63160afbed
3	5	[web, email, mobile]	5	7	bogo
4	5	[web, email]	20	10	discount

```
In [6]: profile.head()

Out[6]:
```

4	68617/ca624614fbc85e91a2a49b52598	offer received	{offer id: 4d5c5/ea9a6940dd891ad53e9d8e8da0}	0
---	-----------------------------------	----------------	----------------------------------------------	---

## Cleaning the portfolio dataframe

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

Out[7]:
```

	person	event	value	time
0	78fa995795e4d85b5d9ceec4395fe	offer received	0	0
1	0a23232639434f34b24c4d3d47e6ba43	offer viewed	0	0
2	e21275569f464592b11af22dc27a7932	offer received	0	0
3	8e9d8bc7a33c4b65b9aef6a799e6d9	offer viewed	0	0
4	68b6c3a43864d31939f3a4f3e3d783	offer received	0	0

## Cleaning the portfolio dataframe

```
In [8]: # Changing the duration column values from days to hours
new_portfolio = portfolio.copy()
new_portfolio['duration'] = new_portfolio['duration'] * 24
```

```
# Applying one hot encoding to the channels column
new_portfolio['web'] = new_portfolio['channels'].apply(lambda x: 1 if 'web' in x else 0)
new_portfolio['email'] = new_portfolio['channels'].apply(lambda x: 1 if 'email' in x else 0)
new_portfolio['mobile'] = new_portfolio['channels'].apply(lambda x: 1 if 'mobile' in x else 0)
new_portfolio['social'] = new_portfolio['channels'].apply(lambda x: 1 if 'social' in x else 0)
```

```
# apply one hot encoding to offer_type column
offer_type = pd.get_dummies(new_portfolio['offer_type'])

# drop the channels and offer_type column
new_portfolio.drop(['channels'], axis=1, inplace=True)
```

```
# combine the portfolio and offer_type dataframe to form a cleaned dataframe
new_portfolio = pd.concat((new_portfolio, offer_type), axis=1, sort=False)
new_portfolio.head()
```

reward	difficulty	duration	offer_type	id	web	email	mobile	social	bogo	discount	informational
0	10	10	168	bogo	a625a4363720a4b09b5b0b210d5d	0	1	1	1	1	0
1	10	10	120	bogo	4d5c7fa8a9404d891ad353e9b8e3a0	1	1	1	1	1	0
2	0	0	96	informational	3f207d87b143eaa3ce63160afbed	1	1	1	0	0	1
3	5	5	168	bogo	9e988bc7a33c4b65b9aef6a799e6d9	1	1	1	0	1	0
4	5	20	240	discount	0b1e15392cc4b7b9fa7c272a2e1d7	1	1	0	0	0	1

## Cleaning the profile dataframe

```
In [9]: #check for null values
profile.isnull().sum()

Out[9]:
```

```
gender      2175
age          0
id           0
became_member_on  2175
income      2175
dtype: int64
```

```
In [10]: profile[profile['age']== 118].age.count()

Out[10]: 2175
```

According to the description of the profile data frame and checking null values, it looks like values of gender & income are missing where age is 118. To confirm, I'll print the values with age 118.

```
In [11]: profile[profile['age']== 118].drop('became_member_on', 'id', axis=1)

Out[11]:
```

gender	age	income
0	None	NaN
2	None	NaN
4	None	NaN
6	None	NaN
7	None	NaN

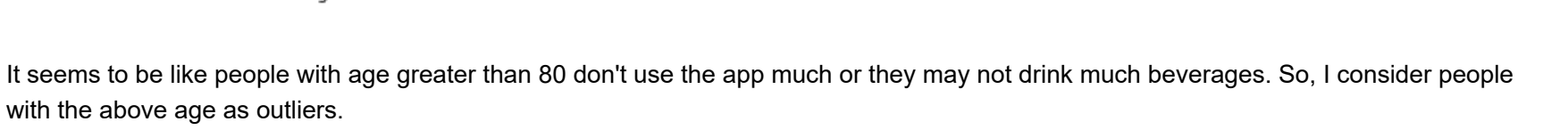
```
-----
0      None    118    NaN
2      None    118    NaN
4      None    118    NaN
6      None    118    NaN
7      None    118    NaN
```

2175 rows x 3 columns

Thus, it is confirmed that both gender & income for the age 118 are missing. All missing age values are encoded as 118

```
In [12]: sns.boxplot(profile['age'], width=0.5);

Out[12]:
```



It seems to be like people with age greater than 80 don't use the app much or they may not drink much beverages. So, I consider people with the above age as outliers.

## Merging the transcript and profile datasets

```
In [13]: #merge the transcript and the profile datasets
df = transcript.merge(profile, left_on='person', right_on='id', how='outer')
df.head()
```

	person	event	value	time	gender	age	ic
0	78fa995795e4d85b5d9ceec4395fe	offer received	0	F	75	78fa995795e4d85b5d9ceec4395fe	
1	78fa995795e4d85b5d9ceec4395fe	offer viewed	6	F	75	78fa995795e4d85b5d9ceec4395fe	
2	78fa995795e4d85b5d9ceec4395fe	transaction	132	F	75	78fa995795e4d85b5d9ceec4395fe	
3	78fa995795e4d85b5d9ceec4395fe	offer completed	132	F	75	78fa995795e4d85b5d9ceec4395fe	
4	78fa995795e4d85b5d9ceec4395fe	transaction	144	F	75	78fa995795e4d85b5d9ceec4395fe	

```
In [14]: null_profile = df[df['gender'].isnull()]
len(null_profile)

Out[14]: 33772
```

```
In [15]: #null profile offer completion rate
null_rate = len(null_profile[null_profile['event']=='offer completed'])/len(null_profile[null_profile['event']=='offer received'])
print(null_rate)

0.1161065466448445
```

Now let us find out what is the offer completion rate for non null profiles.

```
In [16]: not_null_profile = df[df['gender'].notnull()]
not_null_rate = len(not_null_profile[not_null_profile['event']=='offer completed'])/len(not_null_profile[not_null_profile['event']=='offer received'])
print(not_null_rate)

0.4878723628216117
```

The profiles which have null values in their profile have less offer completion rate so, I will be dropping the null values for my analysis.

```
In [17]: #dropping null values
new_df = df[df['gender'].notnull()].copy()
new_df.drop('id', axis=1, inplace=True)
```

```
In [18]: #changing became_member_on to a datetime format
new_df['became_member_on'] = pd.to_datetime(new_df['became_member_on'], format='%Y%m%d')

#creating a new column that has the year which customer became members
new_df['year'] = new_df['became_member_on'].apply(lambda x: str(x)[4:])
```

```
#rounding the time of hours to days and also rounding up
new_df['days'] = new_df['time'].apply(lambda x: int(x / 24) + (x % 24 * 0))

#creating the value column
new_df['value'] = new_df['value'].apply(lambda x: x['offer_id'] if 'offer_id' in x else x['offer_id'] if 'offer_id' in x else np.nan)
new_df['amount'] = new_df['value'].apply(lambda x: x.get('amount', 0))
new_df.drop(['value'], axis=1, inplace=True)
```

```
new_df = new_df.reset_index(drop=True)

In [19]: new_df.head()

Out[19]:
```

	person	event	time	gender	age	became_member_on	income	year	days
0	78fa995795e4d85b5d9ceec4395fe	offer received	0	F	75	2017-05-09	10000.0	2017	0
1	78fa995795e4d85b5d9ceec4395fe	offer viewed	6	F	75	2017-05-09	10000.0	2017	1
2	78fa995795e4d85b5d9ceec4395fe	transaction	132	F	75	2017-05-09	10000.0	2017	6
3	78fa995795e4d85b5d9ceec4395fe	offer completed	132	F	75	2017-05-09	10000.0	2017	6
4	78fa995795e4d85b5d9ceec4395fe	transaction	144	F	75	2017-05-09	10000.0	2017	6

```
In [20]: new_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 272762 entries, 0 to 272761
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 person 272762 non-null object
1 event 272762 non-null object
2 time 272762 non-null int64
3 gender 272762 non-null object
4 age 272762 non-null int64
5 became_member_on 272762 non-null datetime64[ns]
6 income 272762 non-null float64
7 year 272762 non-null object
8 days 272762 non-null int64
9 offer_id 148805 non-null object
10 amount 272762 non-null float64
11 days 272762 non-null int64(3), object(5)
memory usage: 22.9+ MB

new_df.isnull().head()

Out[21]:
```

reward	difficulty	duration	offer_type	id	web	email	mobile	social	bogo	discount	informational
0	10	10	168	bogo	a625a4363720a4b09b5b0b210d5d	0	1	1	1	1	0
1	10	10	120	bogo	4d5c7fa8a9404d891ad353e9b8e3a0	1	1	1	1	1	0
2	0	0	96	informational	3f207d87b143eaa3ce63160afbed	1	1	1	0	0	1
3	5	5	168	bogo	9e988bc7a33c4b65b9aef6a799e6d9	1	1	1	0	1	0
4	5	20	240	discount	0b1e15392cc4b7b9fa7c272a2e1d7	1	1	0	0	0	1

```
In [22]: new_portfolio.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 person 10 non-null object
1 event 10 non-null object
2 time 10 non-null int64
3 gender 10 non-null object
4 age 10 non-null int64
5 became_member_on 10 non-null datetime64[ns]
6 income 10 non-null float64
7 year 10 non-null object
8 days 10 non-null int64
9 offer_id 10 non-null object
10 amount 10 non-null float64
11 days 10 non-null int64(3), object(5)
memory usage: 678.0+ bytes

new_portfolio.isnull().head()

Out[21]:
```

reward	difficulty	duration	offer_type	id	web	email	mobile	social	bogo	discount	informational
0	10	10	168	bogo	a625a4363720a4b09b5b0b210d5d	0	1	1	1	1	0
1	10	10	120	bogo	4d5c7fa8a9404d891ad353e9b8e3a0	1	1	1	1	1	0
2	0	0	96	informational	3f207d87b143eaa3ce63160afbed	1	1	1	0	0	1
3	5	5	168	bogo	9e988bc7a33c4b65b9aef6a799e6d9	1	1	1	0	1	0
4	5	20	240	discount	0b1e15392cc4b7b9fa7c272a2e1d7	1	1	0	0	0	1

```
In [22]: new_portfolio.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 person 10 non-null object
1 event 10 non-null object
2 time 10 non-null int64
3 gender 10 non-null object
4 age 10 non-null int64
5 became_member_on 10 non-null datetime64[ns]
6 income 10 non-null float64
7 year 10 non-null object
8 days 10 non-null int64
9 offer_id 10 non-null object
10 amount 10 non-null float64
11 days 10 non-null int64(3), object(5)
memory usage: 678.0+ bytes

new_portfolio.isnull().head()

Out[21]:
```

reward	difficulty	duration	offer_type	id	web	email	mobile	social	bogo	discount	informational
0	10	10	168	bogo	a625a4363720a4b09b5b0b210d5d	0	1	1	1	1	0
1	10	10	120	bogo	4d5c7fa8a9404d891ad353e9b8e3a0	1	1	1	1	1	0
2	0	0	96	informational	3f207d87b143eaa3ce63160afbed	1	1	1	0	0	1
3	5	5	168	bogo	9e988bc7a33c4b65b9aef6a799e6d9	1	1	1	0	1	0
4	5	20	240	discount	0b1e15392cc4b7b9fa7c272a2e1d7	1	1	0	0	0	1

```
In [22]: new_portfolio.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 person 10 non-null object
1 event 10 non-null object
2 time 10 non-null int64
3 gender 10 non-null object
4 age 10 non-null int64
5 became_member_on 10 non-null datetime64[ns]
6 income 10 non-null float64
7 year 10 non-null object
8 days 10 non-null int64
9 offer_id 10 non-null object
10 amount 10 non-null float64
11 days 10 non-null int64(3), object(5)
memory usage: 678.0+ bytes

new_portfolio.isnull().head()

Out[21]:
```

reward	difficulty	duration	offer_type	id	web	email	mobile	social	bogo	discount	informational
0	10	10	168	bogo	a625a4363720a4b09b5b0b210d5d	0	1	1	1	1	0
1	10	10	120	bogo	4d5c7fa8a9404d891ad353e9b8e3a0	1	1	1	1	1	0
2	0	0	96	informational	3f207d87b143eaa3ce63160afbed	1	1	1	0	0	1
3	5	5	168	bogo	9e988bc7a33c4b65b9aef6a799e6d9	1	1	1	0	1	0
4	5	20	240	discount	0b1e15392cc4b7b9fa7c272a2e1d7	1	1	0	0	0	1

```
In [22]: new_portfolio.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 person 10 non-null object
1 event 10 non-null object
2 time 10 non-null int64
3 gender 10 non-null object
4 age 10 non-null int64
5 became_member_on 10 non-null datetime64[ns]
6 income 10 non-null float64
7 year 10 non-null object
8 days 10 non-null int64
9 offer_id 10 non-null object
10 amount 10 non-null float64
11 days 10 non-null int64(3), object(5)
memory usage: 678.0+ bytes

new_portfolio.isnull().head()

Out[21]:
```

reward	difficulty	duration	offer_type	id	web	email	mobile	social	bogo	discount	informational
0	10	10	168	bogo	a625a4363720a4b09b5b0b210d5d	0	1	1	1	1	0
1	10	10	120	bogo	4d5c7fa8a9404d891ad353e9b8e3a0	1	1	1	1	1	0
2	0	0	96	informational	3f207d87b143eaa3ce63160afbed	1	1	1	0	0	1
3	5	5	168	bogo	9e988bc7a33c4b65b9aef6a799e6d9	1	1	1	0	1	0
4	5	20	240	discount	0b1e15392cc4b7b9fa7c272a2e1d7	1	1	0	0	0	1

```
In [22]: new_portfolio.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 person 10 non-null object
1 event 10 non-null object
2 time 10 non-null int64
3 gender 10 non-null object
4 age 10 non-null int64
5 became_member_on 10 non-null datetime64[ns]
6 income 10 non-null float64
7 year 10 non-null object
8 days 10 non-null int64
9 offer_id 10 non-null object
10 amount 10 non-null float64
11 days 10 non-null int64(3), object(5)
memory usage: 678.0+ bytes

new_portfolio.isnull().head()

Out[21]:
```

reward	difficulty	duration	offer_type	id	web	email	mobile	social	bogo	discount	informational
0	10	10	168	bogo	a625a4363720a4b09b5b0b210d5d	0	1	1	1	1	0
1	10	10	120	bogo	4d5c7fa8a9404d891ad353e9b8e3a0	1	1	1	1	1	0
2	0	0	96	informational	3f207d87b143eaa3ce63160afbed	1	1	1	0	0	1
3	5	5	168	bogo	9e988bc7a33c4b65b9aef6a799e6d9	1	1	1	0	1	0
4	5	20	240	discount	0b1e15392cc4b7b9fa7c272a2e1d7	1	1	0	0	0	1

```
In [22]: new_portfolio.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 person 10 non-null object
1 event 10 non-null object
2 time 10 non-null int64
3 gender 10 non-null object
4 age 10 non-null int64
5 became_member_on 10 non-null datetime64[ns]
6 income 10 non-null float64
7 year 10 non-null object
8 days 10 non-null int64
9 offer_id 10 non-null object
10 amount 10 non-null float64
11 days 10 non-null int64(3), object(5)
memory usage: 678.0+ bytes

new_portfolio.isnull().head()

Out[21]:
```

5	web	10	non-null	int64
6	email	10	non-null	int64
7	mobile	10	non-null	int64
8	social	10	non-null	int64
9	bogo	10	non-null	int8



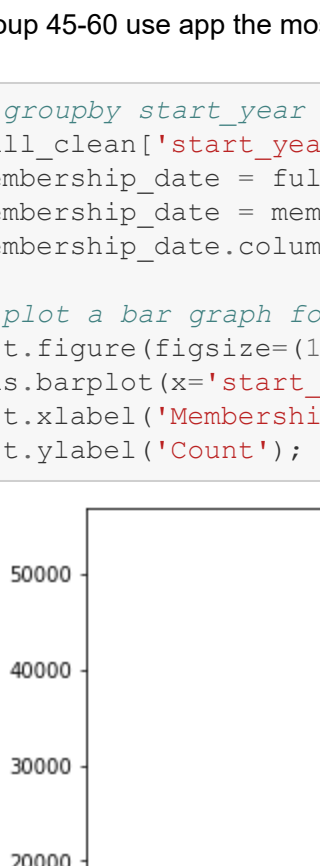
```
In [43]: male_bogo_viewed=full_clean[(full_clean['bogo']==1.0)&(full_clean['event']=='offer viewed')&(full_clean['gender_id']==1)][['event']].count()
female_bogo_viewed=full_clean[(full_clean['bogo']==1.0)&(full_clean['event']=='offer viewed')&(full_clean['gender_id']==0)][['event']].count()
male_bogo_received=full_clean[(full_clean['bogo']==1.0)&(full_clean['event']=='offer received')&(full_clean['gender_id']==1)][['event']].count()
female_bogo_received=full_clean[(full_clean['bogo']==1.0)&(full_clean['event']=='offer received')&(full_clean['gender_id']==0)][['event']].count()
female_success_bogo=((female_bogo_viewed)/(female_bogo_received))*100
male_success_bogo=((male_bogo_viewed)/(male_bogo_received))*100

print('Female: ',female_success_bogo,' Male: ',male_success_bogo)

plt.plot([female_success_bogo,male_success_bogo],colors='lightblue','pink', labels=['males', 'females'], counter-clock=False, shadow=True)
plt.title('Gender likely to view a discount offer')

Female: 83.30751708428245, Male: 82.72619673855866

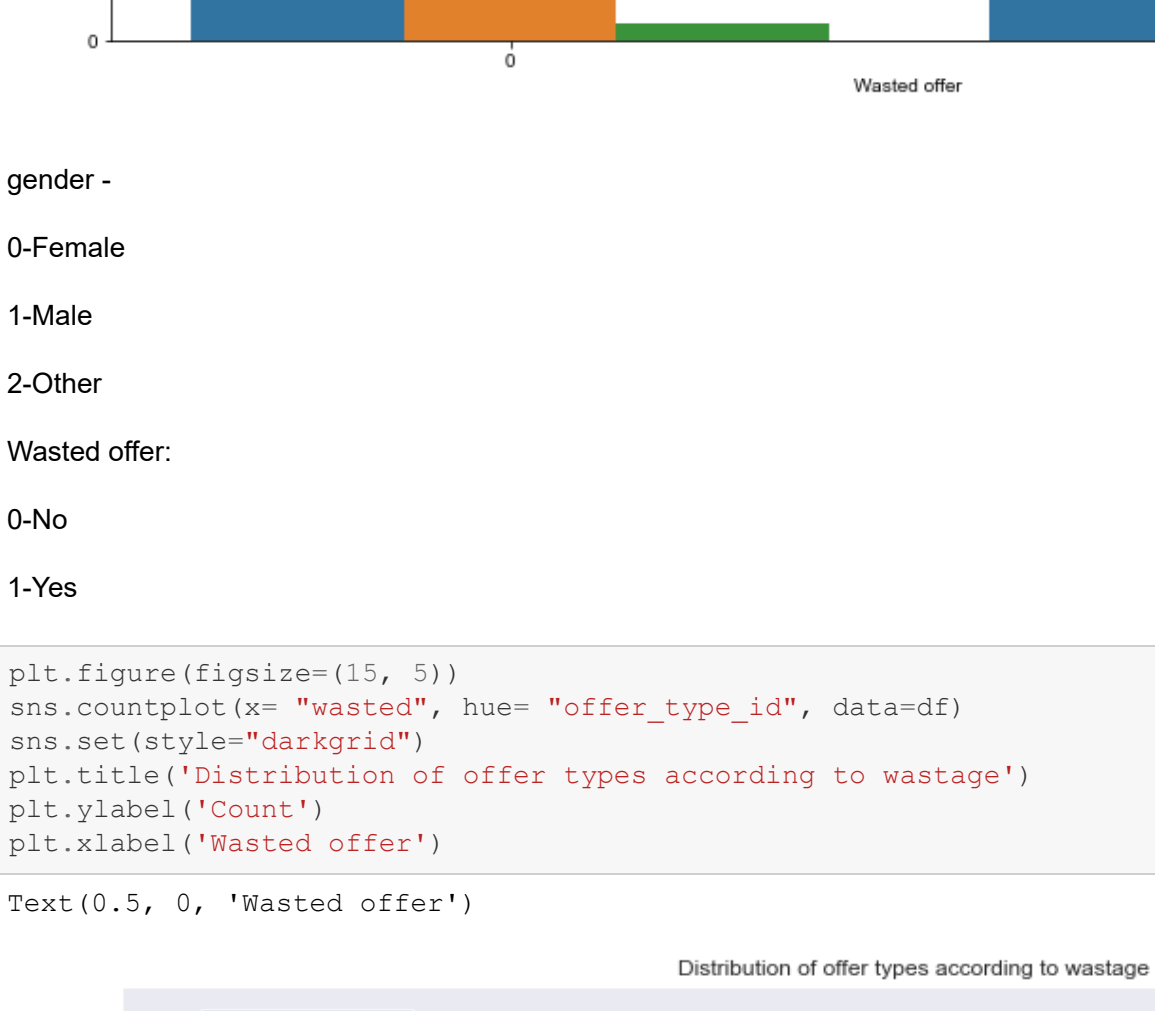
Out[43]: Text(0.5, 1.0, 'Gender likely to view a discount offer')
```



The males represent 62.7% of the data and use the Starbucks app more than the females. Specifically, both males & females in the age group 45-60 use app the most. Discount offers are more preferred by the customers.

```
In [44]: # groupby start year and gender to plot a graph
full_clean['start_year'] = full_clean.became_member_on.dt.year
membership_date = full_clean.groupby(['start_year', 'gender']).size()
membership_date.reset_index(inplace=True)
membership_date.columns = ['start_year', 'gender', 'count']

# plot a bar graph for age distribution as a function of gender in membership program
sns.barplot(x='start_year', y='count', hue='gender', data=membership_date)
plt.xlabel('Membership Start Year')
plt.ylabel('Count')
```



```
In [45]: plt.figure(figsize=(15, 5))
sns.countplot(x='wasted', hue='gender_id', data=df)
sns.set(style='darkgrid')
plt.title('Gender distribution in offer wastage')
plt.ylabel('Count')
plt.xlabel('Wasted offer')
plt.legend(title='Gender')
```



```
In [46]: plt.figure(figsize=(15, 5))
sns.countplot(x='wasted', hue='offer_type_id', data=df)
sns.set(style='darkgrid')
plt.title('Distribution of offer types according to wastage')
plt.ylabel('Count')
plt.xlabel('Wasted offer')
```



```
In [47]: plt.figure(figsize=(15, 5))
sns.countplot(x='wasted', hue='year_id', data=df)
sns.set(style='darkgrid')
plt.title('Year of membership of customers and wastage of offers')
plt.ylabel('Count')
plt.xlabel('Wasted offer')
```

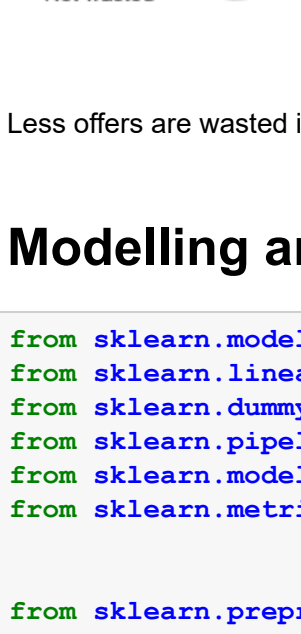


```
In [48]: plt.figure(figsize=(15, 5))
sns.countplot(x='wasted', hue='duration', data=df)
sns.set(style='darkgrid')
plt.title('Duration(in hours) for the offer to be open and wastage of offers')
plt.ylabel('Count')
plt.xlabel('Wasted offer')
```

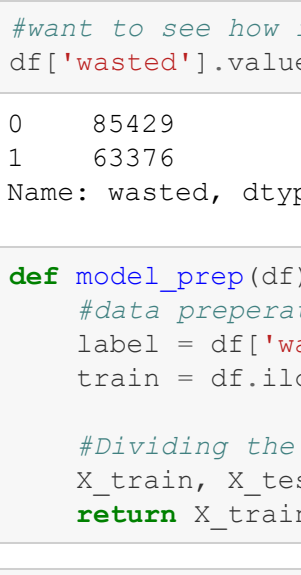


The ideal time for the offer to be open should be 168 hours ie. 7 days for Less number of offers to be wasted

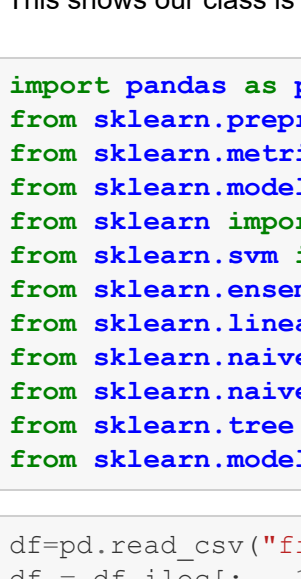
```
In [49]: web_not_wasted=df[(df['wasted']==0)&(df['web']==1.0)][['web']].count()
web_wasted=df[(df['wasted']==1)&(df['web']==1.0)][['web']].count()
plt.pie(web_not_wasted,web_wasted,colors='blue','green', labels=['Not wasted', 'Wasted'], counter-clock=False, shadow=True)
plt.title('Web offer Wastage ')
```



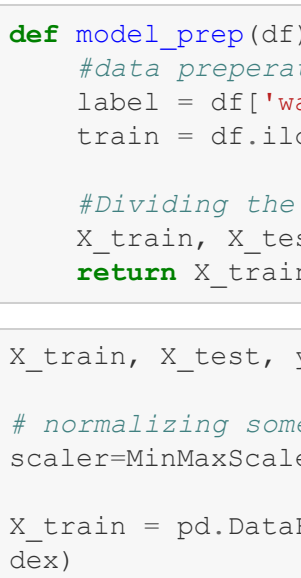
```
In [50]: email_not_wasted=df[(df['wasted']==0)&(df['email']==1.0)][['web']].count()
email_wasted=df[(df['wasted']==1)&(df['email']==1.0)][['web']].count()
plt.pie(email_not_wasted,email_wasted,colors='cyan','pink', labels=['Not wasted', 'Wasted'], counter-clock=False, shadow=True)
plt.title('Email offer wastage ')
```



```
In [51]: mobile_not_wasted=df[(df['wasted']==0)&(df['mobile']==1.0)][['web']].count()
mobile_wasted=df[(df['wasted']==1)&(df['mobile']==1.0)][['web']].count()
plt.pie(mobile_not_wasted,mobile_wasted,colors='pink','yellow', labels=['Not wasted', 'Wasted'], counter-clock=False, shadow=True)
plt.title('Mobile offer wastage ')
```



```
In [52]: social_not_wasted=df[(df['wasted']==0)&(df['social']==1.0)][['web']].count()
social_wasted=df[(df['wasted']==1)&(df['social']==1.0)][['web']].count()
plt.pie(social_not_wasted,social_wasted,colors='black','aqua', labels=['Not wasted', 'Wasted'], counter-clock=False, shadow=True)
plt.title('Social offer wastage ')
```



Less offers are wasted if sent through mobile or web

## Modelling and Evaluation

```
In [53]: from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.linear_model import LogisticRegression
from sklearn.dummy import DummyClassifier
from sklearn.pipeline import Pipeline
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split, StratifiedKFold
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier
from sklearn.metrics import accuracy_score, recall_score, confusion_matrix
import time
```

```
In [54]: #want to see how imbalanced is the dataset
df['wasted'].value_counts()

Out[54]: 0    85429
        1    63376
        Name: wasted, dtype: int64
```

```
In [55]: def model_prep(df):
#data preparation
label = df['wasted']
train = df.iloc[:, 2:].copy()

#Dividing the data into train and test
X_train, X_test, y_train, y_test = train_test_split(train, label, test_size=0.30, random_state=1)
return X_train, X_test, y_train, y_test
```

```
In [56]: X_train, X_test, y_train, y_test = model_prep(df)
#dummy model
dummy1 = DummyClassifier(random_state=1).fit(X_train, y_train)
pred_dummy1 = dummy1.predict(X_test)
print("randomly guessing score: {:.2f}".format(dummy1.score(X_test, y_test)))
dummy2 = DummyClassifier(strategy='most_frequent').fit(X_train, y_train)
pred_dummy2 = dummy2.predict(X_test)
print("Guess all customers will stay score: {:.2f}".format(dummy2.score(X_test, y_test)))

randomly guessing score: 0.51
guess all customers will stay score: 0.57
```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\dummy.py:131: FutureWarning: The default value of strategy will change from stratified to prior in 0.24.
warnings.warn("The default value of strategy will change from "

```
In [1]: import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn import svm
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
```

```
In [2]: df=pd.read_csv('file1.csv')
df = df.iloc[:, 1:]
df.head()
```

age greater than 80 don't use the app much or they may not drink many beverages. There is one entry for age 118. I will be considering it an outlier.

Buy one get one free offer is less likely to get wasted and most of the informational offers are most likely to get wasted. Customers who joined in 2017 are less likely to waste an offer.

The ideal time for the offer to be open should be 168 hours ie. 7 days Less number of offers is wasted if offers are sent through social media or web.

The most common offer type among all age groups is the BOGO , followed by the Discount Offers. Whereas, the least common offer to be sent is the informational offers. I believe that BOGO offers are more attractive compared to other offers provided by Starbucks.

## Model preparation

```
In [4]: def model_prep(df):
#data preparation
label = df['wasted']
train = df.iloc[:, 2:].copy()

#Dividing the data into train and test
X_train, X_test, y_train, y_test = train_test_split(train, label, test_size=0.30, random_state=1)
return X_train, X_test, y_train, y_test
```

```
In [5]: X_train, X_test, y_train, y_test = model_prep(df)
# normalizing some numerical values
scaler=MinMaxScaler()
X_train = pd.DataFrame(MinMaxScaler().fit_transform(X_train), columns=X_train.columns, index=X_train.index)
X_test = pd.DataFrame(MinMaxScaler().fit_transform(X_test), columns=X_test.columns, index=X_test.index)
```

```
In [ ]: model_params = {
'svm': {
'model': svm.SVC(gamma='auto'),
'params': {
'C': (1,10,20),
'kernel': ['rbf','linear']
}
},
'random_forest': {
'model': RandomForestClassifier(),
'params': {
'n_estimators': [1,5,10]
}
},
'logistic_regression': {
'model': LogisticRegression(solver='liblinear', multi_class='auto'),
'params': {
'C': (1,5,10)
}
},
'naive_bayes_gaussian': {
'model': GaussianNB(),
'params': {}
},
'naive_bayes_multinomial': {
'model': MultinomialNB(),
'params': {}
},
'decision_tree': {
'model': DecisionTreeClassifier(),
'params': {
'criterion': ['gini','entropy'],
}
}
}
```

```
In [ ]: #Time4 hrs (app)
scores = []
for model_name, mp in model_params.items():
clf = GridSearchCV(mp['model'], mp['params'], cv=5, return_train_score=False)
clf.fit(X_train, y_train)
scores.append((
'model': model_name,
'best_score': clf.best_score_,
'best_params': clf.best_params_
))

df = pd.DataFrame(scores, columns=['model', 'best_score', 'best_params'])
df.to_csv('file2.csv')
```

Out[8]:

	model	best_score	best_params
0	svm	0.866556	{'C': 20, 'kernel': 'rbf'}
1	random_forest	0.859795	{'n_estimators': 10}
2	logistic_regression	0.826592	{'C': 5}
3	naive_bayes_gaussian	0.806707	{}
4	naive_bayes_multinomial	0.806659	{}
5	decision_tree	0.878007	{'criterion': 'entropy'}

```
In [9]: import matplotlib.pyplot as plt
import seaborn as sns

# draw lineplot
sns.lineplot(x='model', y='best_score', data=df).set_title('Model vs Best scores')
sns.set(rc={'figure.figsize': (12,5)})
plt.show()
```



## Conclusion

The problem that I chose to solve was to build a model that predicts whether a customer will respond to an offer.

The following observations were made after the exploratory data analysis:

The age group 45-60 is the most common in customers. Most Customers have income ranging between 50000-70000. Discount offer is more popular because not only the absolute number of offer completed is slightly higher than BOGO offer, its overall completed/received rate is also about 7% higher.

Females are more likely to complete the bogo or discount offer and males have a low offer completion rate it seems to be that people with age greater than 80 don't use the app much or they may not drink many beverages. There is one entry for age 118. I will be considering it an outlier.

Buy one get one free offer is less likely to get wasted and most of the informational offers are most likely to get wasted. Customers who joined in 2017 are less likely to waste an offer.

The ideal time for the offer to be open should be 168 hours ie. 7 days Less number of offers is wasted if offers are sent through social media or web.

The most common offer type among all age groups is the BOGO , followed by the Discount Offers. Whereas, the least common offer to be sent is the informational offers. I believe that BOGO offers are more attractive compared to other offers provided by Starbucks.

After splitting data to train and test datasets, I chose the best estimator "DecisionTreeClassifier" using GridSearch which is the best performing classifier algorithm among the above 6 classifiers tested. After performing hyper parameter tuning, I built the model. An accuracy of 87.8 was achieved by the model.

```
In [ ]:
```