

# Introduction

Our project is to forecast those customers who will become loyal to the product based on the processed data of offering incentives (a.k.a. coupons) to a large number of customers. The customers who are offered a discount and redeem the offer are the focus of this project. We are going to build models to predict which of the target customers will return to purchase the same item again.

The dataset we have includes a minimum of a year of shopping history prior to each customer's incentive, as well as the purchase histories of many other shoppers. The transaction history contains all items purchased, not just items related to the offer. Only one offer per customer is included in the data. The training set is comprised of offers issued before 2013-05-01. The test set is offers issued on or after 2013-05-01.

The workflow of this project includes data preprocessing, data manipulation, feature engineering, exploratory data analysis and model building. We create various predictive models including logistic regression, decision tree, random forests and xgboost. We also obtain the feature importance and performance measurement after obtaining the model results.

## Initial set-up & data preprocessing

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn import preprocessing
4 import matplotlib.pyplot as plt
5 plt.rc("font", size=14)
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.cross_validation import train_test_split
8 from sklearn import metrics
9 from sklearn import datasets
10 from sklearn.feature_selection import RFE
11 #import seaborn as sns
12 #sns.set(style="white")
13 #sns.set(style="whitegrid", color_codes=True)
```

```
/anaconda3/lib/python3.6/site-packages/sklearn/cross_validation.py:41:
DeprecationWarning: This module was deprecated in version 0.18 in favo
r of the model_selection module into which all the refactored classes
and functions are moved. Also note that the interface of the new CV it
erators are different from that of this module. This module will be re
moved in 0.20.
```

```
"This module will be removed in 0.20.", DeprecationWarning)
```

## Reduce the size of the data

Here we select chain 2, 4 and 8. Save the file as reduced2.csv

```
In [ ]:

1  #Chosing transactions from 3 chains of the store.
2  loc_transactions = "data/transactions.csv"
3  loc_reduced = "data/reduced2.csv" #Output file
4  def reduce_data(loc_transactions, loc_reduced):
5      with open(loc_reduced, "w") as outfile:
6          for e, line in enumerate( open(loc_transactions) ):
7              if line.split(",")[1] in ['2','4','8']: # Add more chains by adding
8                  outfile.write( line )
9
10 reduce_data(loc_transactions, loc_reduced)
11 # reduced2.csv contains transactions from the chain ids mentioned in list above
12 # redcued2.csv is your subset of the transactions data. Look at the customers p
```

## Data Loading

1. reduced transaction data - here we selected only 3 chains with id 2, 4 and 8.

```
In [35]:

1  data=pd.read_csv("data/reduced2.csv",header=None,names=["id","chain","dept","cat
2  data.head()
3
```

Out[35]:

	id	chain	dept	category	company	brand	date	productsize	productmeasure	pro
0	12524696	4	58	5833	103760030	8501	2012-03-02	24.00	OZ	
1	12524696	4	58	5833	103760030	8501	2012-03-02	18.40	OZ	
2	12524696	4	99	9901	107989373	12908	2012-03-02	0.66	OZ	
3	12524696	4	57	5710	103663232	4568	2012-03-06	5.30	OZ	
4	12524696	4	9	907	102113020	15704	2012-03-06	16.00	OZ	

2. Training data - containing incentives offered to each customer and their response afterwards

In [36]:

```
1 #Read the training data
2 #train=pd.read_csv("trainHistory.csv",header=1,names=["id","chain","offer",'mar
3 train = pd.read_csv("data/trainHistory.csv")
4 print(len(train))
5 train.head()
```

160057

Out[36]:

	id	chain	offer	market	repeattrips	repeater	offerdate
0	86246	205	1208251	34	5	t	2013-04-24
1	86252	205	1197502	34	16	t	2013-03-27
2	12682470	18	1197502	11	0	f	2013-03-28
3	12996040	15	1197502	9	0	f	2013-03-25
4	13089312	15	1204821	9	0	f	2013-04-01

In [5]:

```
1 len(train['id'].unique())
```

Out[5]:

160057

### 3. Offer data

In [38]:

```
1 #Read the offers data
2 offer=pd.read_csv("data/offers.csv",header=1,names=['offer','category','quantity
3 offer.head()
```

Out[38]:

	offer	category	quantity	company	offervalue	brand
0	1194044	9909	1	107127979	1.00	6732
1	1197502	3203	1	106414464	0.75	13474
2	1198271	5558	1	107120272	1.50	5072
3	1198272	5558	1	107120272	1.50	5072
4	1198273	5558	1	107120272	1.50	5072

# Data preprocessing

In [39]:

```
1 #Count how many people will return to store in training data
2 train['repeater'].value_counts()
```

Out[39]:

```
f      116619
t       43438
Name: repeater, dtype: int64
```

In [40]:

```
1 (train['repeater'][1]=='f')
```

Out[40]:

False

In [41]:

```
1 #Convert the value in repeater column into integer
2 train.repeater = (train.repeater == 't').astype(int)
3 train.head()
```

Out[41]:

	id	chain	offer	market	repeattrips	repeater	offerdate
0	86246	205	1208251	34	5	1	2013-04-24
1	86252	205	1197502	34	16	1	2013-03-27
2	12682470	18	1197502	11	0	0	2013-03-28
3	12996040	15	1197502	9	0	0	2013-03-25
4	13089312	15	1204821	9	0	0	2013-04-01

In [11]:

```
1 print(train.shape)
2 print(data.shape)
3 print(len(set(train.id)&set(data.id)))
```

```
(160057, 7)
(9480946, 11)
4732
```

Merge Offer data with Training data on "offer" and save the merged data as train\_offer.

In [42]:

```
1 train_offer = pd.merge(train,offer,how='left',on='offer')
2 print(len(train_offer[train_offer['company'].notnull()]))
3 train_offer = train_offer.rename(columns={'company': 'offer_company','brand':'offer_brand'})
4 train_offer.head()
```

160057

Out[42]:

	id	chain	offer	market	repeattrips	repeater	offerdate	offer_category	quantity	c
0	86246	205	1208251	34	5	1	2013-04-24	2202	1	
1	86252	205	1197502	34	16	1	2013-03-27	3203	1	
2	12682470	18	1197502	11	0	0	2013-03-28	3203	1	
3	12996040	15	1197502	9	0	0	2013-03-25	3203	1	
4	13089312	15	1204821	9	0	0	2013-04-01	5619	1	

In [10]:

```
1 print(train_offer.shape)
2 print(train.shape)
3 print(len(train.id.unique()))
4 print(len(data.id.unique()))
```

(160057, 12)

(160057, 7)

160057

8946

In [13]:

```
1 print('transaction:',data.columns)
2 print('train: ',train.columns)
3 print('offer: ',offer.columns)
4 print('offer: ',train_offer.columns)
```

```
transaction: Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'date',
                    'productsize', 'productmeasure', 'productquantity', 'productamount'],
                  dtype='object')
train: Index(['id', 'chain', 'offer', 'market', 'repeattrips', 'repeater',
             'offerdate'],
            dtype='object')
offer: Index(['offer', 'category', 'quantity', 'company', 'offervalue', 'brand'], dtype='object')
offer: Index(['id', 'chain', 'offer', 'market', 'repeattrips', 'repeater',
             'offerdate', 'offer_category', 'quantity', 'offer_company',
             'offervalue', 'offer_brand'],
            dtype='object')
```

In [121]:

```
1 def univariate_plotter(feature='No feature passed', train_sub_col_target=0, target_col=1):
2     X_train = train_sub_col_target.copy()
3     X_train['y'] = train_sub_col_target[target_col]
4     #for col_number in range(7, 8):
5     col = feature #X_train.columns[col_number]
6     print('Below are the actual found rates wrt ' + col)
7     nan_flag=0
8     if pd.isnull(X_train[col]).sum()>0:
9         nan_flag=1
10        print("NANs present")
11    bins = bins
12    # cuts=[0]
13    cuts=[]
14    prev_cut = -1000000000
15    if nan_flag!=1:
16        for i in range(bins+1):
17            next_cut = np.percentile(train_sub_col_target[col], i*100/bins)
18            if next_cut!=prev_cut:
19                cuts.append(next_cut)
20            prev_cut = next_cut#.copy()
21        #print(cuts)
22        cuts[0] = cuts[0]-1
23        cuts[len(cuts)-1] = cuts[len(cuts)-1]+1
24        #cuts=[ -1, 0.01, 0.015, 0.02, 0.025, 0.03, 0.035, 0.04, 0.05, 2]
25        cut_series = pd.cut(X_train[col], cuts)
26    else:
27        train_sub_col_no_nan = train_sub_col_target[col].copy()
```

```

28 train_sub_col_no_nan[np.isinf(train_sub_col_no_nan)]=np.nan
29 train_sub_col_no_nan = train_sub_col_no_nan[~np.isnan(train_sub_col_no_n
30 X_train[col][np.isinf(X_train[col])]=np.nan
31 X_train_no_nan = X_train[~np.isnan(X_train[col])]
32
33 for i in range(bins+1):
34     next_cut = np.percentile(train_sub_col_no_nan, i*100/bins)
35     if next_cut!=prev_cut:
36         cuts.append(next_cut)
37     prev_cut = next_cut.copy()
38
39 cuts[0] = cuts[0]-0.00001
40 cuts[len(cuts)-1] = cuts[len(cuts)-1]+0.00001
41 #cuts=[ -1, 0.01, 0.015, 0.02, 0.025, 0.03, 0.035, 0.04, 0.05]
42 cut_series = pd.cut(X_train_no_nan[col], cuts)
43 cut_series_test = pd.cut(X_test_no_nan[col], cuts)
44 #print(cuts)
45 if nan_flag!=1:
46
47     grouped = X_train.groupby([cut_series], as_index=True).agg({'y':[np.size
48     grouped1 = pd.DataFrame(grouped.index)
49     grouped = X_train.groupby([cut_series], as_index=False).agg({'y':[np.si
50     grouped.columns = ['_'.join(cols).strip() for cols in grouped.columns.va
51
52     grouped = pd.DataFrame(grouped.to_records())
53     grouped1['y_mean'] = grouped['y_mean']
54     grouped1['y_sum'] = grouped['y_size']
55     grouped1[col+'_mean'] = grouped[col+'_mean']
56
57 else:
58     grouped = X_train_no_nan.groupby([cut_series], as_index=True).agg({'y':
59     grouped1 = pd.DataFrame(grouped.index)
60     grouped = X_train_no_nan.groupby([cut_series], as_index=False).agg({'y'
61     grouped.columns = ['_'.join(cols).strip() for cols in grouped.columns.va
62     grouped = pd.DataFrame(grouped.to_records())
63     grouped1['y_mean'] = grouped['y_mean']
64     grouped1['y_sum'] = grouped['y_sum']
65     grouped1_nan = grouped1[0:1]
66     grouped1_nan[col] = (grouped1_nan[col]).astype('str')
67     grouped1_nan[col][0] = 'Nan'
68     grouped1_nan['y_mean'][0] = y_train[np.isnan(X_train[col])].mean()
69     grouped1_nan['y_sum'][0] = y_train[np.isnan(X_train[col])].sum()
70     grouped1 = pd.concat([grouped1, grouped1_nan])
71
72 grouped1=grouped1.reset_index(drop=True)
73 a = plt.plot(grouped1['y_mean'], marker = 'o')
74 plt.xticks(np.arange(len(grouped1)), (grouped1[col]).astype('str'), rotation
75 plt.xlabel('Bins of '+feature)
76 plt.ylabel('Bin-wise percentage repeater')
77 plt.show()
78 b = plt.bar(np.arange(len(grouped1)), grouped1['y_sum'], alpha=0.5)
79 plt.xticks(np.arange(len(grouped1)), (grouped1[col]).astype('str'), rotation
80 plt.xlabel('Bins of '+feature)
81 plt.ylabel('Bin-wise Population')

```

```
82 plt.show()
83 return(grouped1)
```

Merge train\_offer with Reduced transaction data on customer ID

In [43]:

```
1 reduced_transaction_offer = pd.merge(data,train_offer[['id','offer_company','offerdate'],
2
```

In [30]:

```
1 reduced_transaction_offer.head()
```

Out[30]:

	id	chain	dept	category	company	brand	date	productsize	productmeasure	product
0	12524696	4	58	5833	103760030	8501	2012-03-02	24.00	OZ	
1	12524696	4	58	5833	103760030	8501	2012-03-02	18.40	OZ	
2	12524696	4	99	9901	107989373	12908	2012-03-02	0.66	OZ	
3	12524696	4	57	5710	103663232	4568	2012-03-06	5.30	OZ	
4	12524696	4	9	907	102113020	15704	2012-03-06	16.00	OZ	

In [44]:

```
1 b = reduced_transaction_offer.groupby(['id'],as_index = False).agg({'offerdate':
2 b.head()
```

Out[44]:

	id	offerdate
0	12524696	NaN
1	13501141	NaN
2	13744500	NaN
3	13807224	2013-04-05
4	13873775	2013-03-26



# Feature Engineering

Find the most recent purchase date of each customer

In [45]:

```
1 has_bought_ = reduced_transaction_offer.groupby(['id'],as_index = False).agg({'c
```

In [27]:

```
1 has_bought_.head()
```

Out[27]:

	id	date
0	12524696	2013-06-19
1	13501141	2013-05-08
2	13744500	2013-06-15
3	13807224	2013-04-03
4	13873775	2013-03-23

In [88]:

```
1 print(has_bought_.shape)
2 print(has_bought_[has_bought_['date'].notnull()].shape)
```

(8946, 2)
(8946, 2)

Additional feature: Most recent purchasing date

In [46]:

```
1 train_offer_ = pd.merge(train_offer,has_bought_[['id','date']],
2                           how = 'left',on = 'id')
3 print(train_offer_[train_offer_['date'].notnull()].shape)
4 train_offer_ =(train_offer_[train_offer_['date'].notnull()])
5 train_offer_.head()
```

(4732, 13)

Out[46]:

	id	chain	offer	market	repeattrips	repeater	offerdate	offer_category	quantity
8	13807224	4	1204576	1	0	0	2013-04-05	5616	1
9	13873775	4	1197502	1	0	0	2013-03-26	3203	1
10	13974451	4	1197502	1	0	0	2013-03-26	3203	1
12	14381137	4	1197502	1	0	0	2013-04-04	3203	1
19	15994113	4	1197502	1	0	0	2013-03-26	3203	1

**Additional feature: Delta.** delta is the number of 30 days, i.e. if delta = 1, then it means 30 days, if delta = 60, it means 60 days and etc.

In [47]:

```
1 from datetime import datetime
2 import math
3 def __datetime(date_str):
4     return datetime.strptime(date_str, '%Y-%m-%d')
5 #delta is the number of 30 days, i.e. if delta = 1, then it means 30 days, if delta = 60, it means 60 days and etc.
6 delta = (train_offer_.offerdate.apply(__datetime)-train_offer_.date.apply(__datetime)).dt.days//30
7 train_offer_.columns
```

Out[47]:

Index(['id', 'chain', 'offer', 'market', 'repeattrips', 'repeater', 'offerdate', 'offer\_category', 'quantity', 'offer\_company', 'offervalue', 'offer\_brand', 'date'], dtype='object')

In [48]:

```
1 delta.shape
2 delta1 = delta
3 #delta1 = pd.to_timedelta(delta,unit='d').astype('timedelta64[D]')
4 delta1 = delta1/30
5 delta1 = delta1.dt.ceil('D')
6 delta1 = delta1.dt.days
7 #delta1.head()
```

In [92]:

```
1 print(reduced_transaction_offer.shape)
2 print(reduced_transaction_offer.columns)
```

(9480946, 15)  
Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'date',  
 'productsize', 'productmeasure', 'productquantity', 'productamo  
unt',  
 'offer\_company', 'offer\_brand', 'offer\_category', 'offerdate'],  
 dtype='object')

In [49]:

```
1 ##### Find transactional data with company of the transaction equal to company c
2 reduced_transaction_offer_ = reduced_transaction_offer[reduced_transaction_offer
3                                                         reduced_transaction_offer
4 reduced_transaction_offer_.head()
```

Out[49]:

	id	chain	dept	category	company	brand	date	productsize	productmeasure	
5746	13807224	4	56	5610	104610040	15889	2012-06-25	5.00	OZ	
5845	13807224	4	56	5616	104610040	15889	2012-08-01	6.67	OZ	
5846	13807224	4	56	5616	104610040	15889	2012-08-01	7.50	OZ	
5996	13807224	4	56	5616	104610040	15889	2012-09-17	7.60	OZ	
5997	13807224	4	56	5616	104610040	15889	2012-09-17	6.67	OZ	

In [94]:

```
1 print(len(reduced_transaction_offer_.id.unique()))
2 #print(len(reduced_transaction_offer_brand.id.unique()))
```

In [95]:

```
1 print(reduced_transaction_offer_.shape)
2 print(reduced_transaction_offer_.columns)
```

(12637, 15)

Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'date', 'productsize', 'productmeasure', 'productquantity', 'productamount', 'offer\_company', 'offer\_brand', 'offer\_category', 'offerdate'], dtype='object')

**Calculate total amount and number of times each customer has bought from each company on offer**

In [50]:

```
1 #Total amount and number of times each customer spent in each company
2 has_bought_ = reduced_transaction_offer_.groupby(['id'],as_index = False).agg(n
3 has_bought_amount = reduced_transaction_offer_.groupby(['id'],as_index = False)
4 #print(has_bought_company.head())
5 print(has_bought_.columns)
6 #has_bought_company_amount.head()
7 has_bought_.head()
```

Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'date', 'productsize', 'productmeasure', 'productquantity', 'productamount', 'offer\_company', 'offer\_brand', 'offer\_category', 'offerdate'], dtype='object')

Out[50]:

	id	chain	dept	category	company	brand	date	productsize	productmeasure	productquantity
0	13807224	7	7	7	7	7	7	7.0	7	7
1	13873775	16	16	16	16	16	16	16.0	16	16
2	16138642	2	2	2	2	2	2	2.0	2	2
3	16535563	3	3	3	3	3	3	3.0	3	3
4	18277411	15	15	15	15	15	15	15.0	15	15

**Additional features: 1. has\_bought\_company; 2. has\_bought\_company\_amount, 3. has\_bought\_company\_quantity**

In [51]:

```
1 train_offer_ = pd.merge(train_offer_,has_bought_[['id','chain']],
2                           how = 'left',on = 'id')
3 train_offer_ = pd.merge(train_offer_,has_bought_amount[['id','productamount'],'p
4                           how = 'left',on = 'id')
5
6 train_offer_.head()
```

Out[51]:

	id	chain_x	offer	market	repeattrips	repeater	offerdate	offer_category	quantity
0	13807224	4	1204576	1	0	0	2013-04-05	5616	1
1	13873775	4	1197502	1	0	0	2013-03-26	3203	1
2	13974451	4	1197502	1	0	0	2013-03-26	3203	1
3	14381137	4	1197502	1	0	0	2013-04-04	3203	1
4	15994113	4	1197502	1	0	0	2013-03-26	3203	1

In [98]:

```
1 train_offer_.columns
```

Out[98]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'chain_y', 'productamount',
      'productquantity'],
      dtype='object')
```

In [52]:

```
1 print(train_offer_.shape)
2 #print(len([reduced_transaction_offer['chain_y'].notnull()]))
3 #print(reduced_transaction_offer[reduced_transaction_offer['chain_y'].notnull()])
4 train_offer_ = train_offer_.rename(
5     columns={'chain_y': 'has_bought_company', 'productamount': 'has_bought_company'})
6 train_offer_.head()
```

(4732, 16)

Out[52]:

	id	chain_x	offer	market	repeattrips	repeater	offerdate	offer_category	quantity
0	13807224	4	1204576	1	0	0	2013-04-05	5616	1
1	13873775	4	1197502	1	0	0	2013-03-26	3203	1
2	13974451	4	1197502	1	0	0	2013-03-26	3203	1
3	14381137	4	1197502	1	0	0	2013-04-04	3203	1
4	15994113	4	1197502	1	0	0	2013-03-26	3203	1

Find transactional data with brand of the transaction equal to brand of the offer

In [53]:

```
1 reduced_transaction_offer_ = reduced_transaction_offer[reduced_transaction_offer['chain_y'] == reduced_transaction_offer['chain_x']]
2
3 reduced_transaction_offer_.shape
```

Out[53]:

(8858, 15)

Calculate total amount and number of times each customer has bought from each brand on an offer

In [54]:

```
1 has_bought_ = reduced_transaction_offer_.groupby(['id'],as_index = False).agg(n
2 has_bought_amount = reduced_transaction_offer_.groupby(['id'],as_index = False)
3 #print(has_bought_brand.head())
4 print(has_bought_amount.shape)
5 print(has_bought_.columns)
```

```
(1564, 12)
Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'date',
      'productsize', 'productmeasure', 'productquantity', 'productamo
unt',
      'offer_company', 'offer_brand', 'offer_category', 'offerdate'],
      dtype='object')
```

**Additional features: 1. has\_bought\_brand; 2. has\_bought\_brand\_amount, 3. has\_bought\_brand\_quantity**

In [55]:

```
1 train_offer_ = pd.merge(train_offer_,has_bought_[['id','chain']],
2                          how = 'left',on = 'id')
3 train_offer_ = pd.merge(train_offer_,has_bought_amount[['id','productamount','p
4                          how = 'left',on = 'id')
5 train_offer_.columns
```

Out[55]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'chain',
      'productamount', 'productquantity'],
      dtype='object')
```

In [56]:

```
1 print(train_offer_.shape)
2 train_offer_ = train_offer_.rename(
3     columns={'chain': 'has_bought_brand', 'productamount': 'has_bought_brand_amou
4             'productquantity': 'has_bought_brand_q'})
5 train_offer_.head()
6
7 print(train_offer_[train_offer_['has_bought_brand'].notnull()].shape)
```

(4732, 19)

(1564, 19)

In [57]:

```
1 reduced_transaction_offer_ = reduced_transaction_offer[reduced_transaction_offer
2                                     reduced_transaction_offer
3 reduced_transaction_offer_.shape
```

Out[57]:

(16817, 15)

In [58]:

```
1 has_bought_ = reduced_transaction_offer_.groupby(['id'],as_index = False).agg(n
2 has_bought_amount = reduced_transaction_offer_.groupby(['id'],as_index = False)
3 #print(has_bought_brand.head())
4 print(has_bought_.columns)
5
6 print(has_bought_amount.columns)
```

```
Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'date',
      'productsize', 'productmeasure', 'productquantity', 'productamo
unt',
      'offer_company', 'offer_brand', 'offer_category', 'offerdate'],
      dtype='object')
Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'product
size',
      'productquantity', 'productamount', 'offer_company', 'offer_bra
nd',
      'offer_category'],
      dtype='object')
```

```
1 ##### Additional features: 1. has_bought_category; 2. has_bought_category_amount, 3.
has_bought_category_quantity
```

In [59]:

```
1 train_offer_ = pd.merge(train_offer_,has_bought_[['id','chain']],
2                           how = 'left',on = 'id')
3 train_offer_ = pd.merge(train_offer_,has_bought_amount[['id','productamount'],'p
4                           how = 'left',on = 'id')
5
```

In [60]:

```
1 print(train_offer_.shape)
2 train_offer_ = train_offer_.rename(
3     columns={'chain': 'has_bought_category', 'productamount': 'has_bought_category
4 #reduced_transaction_offer.head()
5 print(train_offer_[train_offer_['has_bought_category'].notnull()].shape)
```

(4732, 22)

(1861, 22)



In [61]:

```
1 reduced_transaction_offer_ = reduced_transaction_offer[reduced_transaction_offer
2                                     reduced_transaction_offer
3 reduced_transaction_offer_ = reduced_transaction_offer_[reduced_transaction_offe
4                                     reduced_transaction_offer_
5 reduced_transaction_offer_ = reduced_transaction_offer_[reduced_transaction_offe
6                                     reduced_transaction_offer_
```

In [62]:

```
1 has_bought_ = reduced_transaction_offer_.groupby(['id'],as_index = False).agg(np
2 has_bought_amount = reduced_transaction_offer_.groupby(['id'],as_index = False)
3 #print(has_bought_brand.head())
4 print(has_bought_.columns)
5
6 #print(has_bought_amount.columns)
```

```
Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'date',
      'productsize', 'productmeasure', 'productquantity', 'productamo
unt',
      'offer_company', 'offer_brand', 'offer_category', 'offerdate'],
      dtype='object')
```

In [63]:

```
1 train_offer_ = pd.merge(train_offer_,has_bought_[['id','chain']],
2                           how = 'left',on = 'id')
```

In [64]:

```
1 train_offer_ = pd.merge(train_offer_,has_bought_amount[['id','productamount','p
2                           how = 'left',on = 'id')
3 train_offer_.columns
```

Out[64]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bough
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'chain',
      'productamount', 'productquantity'],
      dtype='object')
```

**Additional feature: 1. has\_bought\_all 2. has\_bought\_all\_amount 3. has\_bought\_all\_q**

In [65]:

```
1 print(train_offer_.shape)
2 train_offer_ = train_offer_.rename(
3     columns={'chain': 'has_bought_all', 'productamount': 'has_bought_all_amount',
4 #reduced_transaction_offer.head()
5
6 print(train_offer_[train_offer_['has_bought_category'].notnull()].shape)
```

(4732, 25)

(1861, 25)

In [66]:

```
1 train_offer_ = train_offer_.assign(Delta = delta1)
```

In [42]:

```
1 train_offer_.columns
```

Out[42]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bough
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta'],
      dtype='object')
```

In [35]:

```
1 train_offer_.Delta.unique()
```

Out[35]:

```
array([ nan,   1.,   2.,   3.] )
```

In [114]:

```
1 print(train_offer_.columns)
2 print(offer.columns)
```

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bought_
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta'],
      dtype='object')
Index(['offer', 'category', 'quantity', 'company', 'offervalue', 'bran
d'], dtype='object')
```

In [67]:

```
1 has_bought_ = reduced_transaction_offer.groupby(['id'],as_index = False).agg(np
2 has_bought_amount = reduced_transaction_offer.groupby(['id'],as_index = False).agg
```

In [37]:

```
1 print(has_bought_.shape)
2 print(has_bought_.columns)
3 print(has_bought_amount.columns)
4 print(reduced_transaction_offer.shape)
5 print(len(reduced_transaction_offer.id.unique()))
6 #has_bought_.head()
```

```
(8946, 15)
Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'date',
      'productsize', 'productmeasure', 'productquantity', 'productamo
unt',
      'offer_company', 'offer_brand', 'offer_category', 'offerdate'],
      dtype='object')
Index(['id', 'chain', 'dept', 'category', 'company', 'brand', 'product
size',
      'productquantity', 'productamount', 'offer_company', 'offer_bra
nd',
      'offer_category'],
      dtype='object')
(9480946, 15)
8946
```

**Additional feature: 1. total\_amount, 2. total\_q**

In [68]:

```
1 train_offer_ = pd.merge(train_offer_,has_bought_[['id','chain']],how='left',on=
2 train_offer_ = pd.merge(train_offer_,has_bought_amount[['id','productquantity',
3 train_offer_.columns
```

Out[68]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bough
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'chain',
      'productquantity', 'productamount'],
      dtype='object')
```

In [69]:

```
1 train_offer_ = train_offer_.rename(
2     columns={'chain': 'total_times', 'productamount': 'total_amount', 'productquant
3 #reduced_transaction_offer.head()
```

In [41]:

```
1 print(train_offer_.shape)
2 train_offer_.head()
```

(4732, 29)

Out[41]:

	id	chain_x	offer	market	repeattrips	repeater	offerdate	offer_category	quantity
0	13807224	4	1204576	1	0	0	2013-04-05	5616	1
1	13873775	4	1197502	1	0	0	2013-03-26	3203	1
2	13974451	4	1197502	1	0	0	2013-03-26	3203	1
3	14381137	4	1197502	1	0	0	2013-04-04	3203	1
4	15994113	4	1197502	1	0	0	2013-03-26	3203	1

5 rows × 29 columns

**Additional features: 1. has\_bought\_company\_30 2. has\_bought\_company\_60 3. has\_bought\_company\_90**

In [70]:

```
1 a = train_offer[train_offer['Delta'] == 1]
2 a = a[a['has_bought_company'] > 0]
3 a['has_bought_company_30'] = 1
4 a.columns
```

Out[70]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bought_
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_ti
mes',
      'total_q', 'total_amount', 'has_bought_company_30'],
      dtype='object')
```

In [71]:

```
1 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_company_30']], how='left')
```

In [72]:

```
1 a = train_offer_[train_offer_['Delta'] == 2]
2 a = a[a['has_bought_company'] > 0]
3 a['has_bought_company_60'] = 1
4 #a.columns
5 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_company_60']], how='left')
6 train_offer_.columns
```

Out[72]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bought_
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_ti
mes',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60'],
      dtype='object')
```

In [73]:

```
1 a = train_offer_[train_offer_['Delta'] == 3]
2 a = a[a['has_bought_company'] > 0]
3 a['has_bought_company_90'] = 1
4 #a.columns
5 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_company_90']], how='left')
6 train_offer_.columns
```

Out[73]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bough
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_ti
mes',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60', 'has_bought_company_90'],
      dtype='object')
```

**Additional features: 1. has\_bought\_brand\_30 2. has\_bought\_brand\_60 3. has\_bought\_brand\_90**

In [74]:

```
1 a = train_offer[train_offer['Delta'] == 1]
2 a = a[a['has_bought_brand'] > 0]
3 a['has_bought_brand_30'] = 1
4 #a.columns
5 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_brand_30']], how='left', on='id')
6 train_offer_.columns
```

Out[74]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bought_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_category',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bought_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_times',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60', 'has_bought_company_90',
      'has_bought_brand_30'],
      dtype='object')
```

In [75]:

```
1 a = train_offer_[train_offer_['Delta'] == 2]
2 a = a[a['has_bought_brand'] > 0]
3 a['has_bought_brand_60'] = 1
4 #a.columns
5 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_brand_60']], how='left', on='id')
6 train_offer_.columns
```

Out[75]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bought_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_category',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bought_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_times',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60', 'has_bought_company_90', 'has_bought_brand_30',
      'has_bought_brand_60'],
      dtype='object')
```



In [76]:

```
1 a = train_offer_[train_offer_['Delta'] == 3]
2 a = a[a['has_bought_brand'] > 0]
3 a['has_bought_brand_90'] = 1
4 #a.columns
5 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_brand_90']], how='left', on='id')
6 train_offer_.columns
```

Out[76]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bought_
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_ti
mes',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60', 'has_bought_company_90', 'has_bought_b
rand_30',
      'has_bought_brand_60', 'has_bought_brand_90'],
      dtype='object')
```

**Additional features: 1. has\_bought\_category\_30 2. has\_bought\_category\_60 3. has\_bought\_category\_90**

In [77]:

```
1 a = train_offer_[train_offer_['Delta'] == 1]
2 a = a[a['has_bought_category'] > 0]
3 a['has_bought_category_30'] = 1
4 #a.columns
5 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_category_30']], how='left'
6 train_offer_.columns
```

Out[77]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bough
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_ti
mes',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60', 'has_bought_company_90', 'has_bought_b
rand_30',
      'has_bought_brand_60', 'has_bought_brand_90', 'has_bought_categ
ory_30'],
      dtype='object')
```

In [78]:

```
1 a = train_offer_[train_offer_['Delta'] == 2]
2 a = a[a['has_bought_category'] > 0]
3 a['has_bought_category_60'] = 1
4 #a.columns
5 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_category_60']], how='left'
6 train_offer_.columns
```

Out[78]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bough
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_ti
mes',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60', 'has_bought_company_90', 'has_bought_b
rand_30',
      'has_bought_brand_60', 'has_bought_brand_90', 'has_bought_categ
ory_30',
      'has_bought_category_60'],
      dtype='object')
```

In [79]:

```
1 a = train_offer_[train_offer_['Delta'] == 3]
2 a = a[a['has_bought_category'] > 0]
3 a['has_bought_category_90'] = 1
4 #a.columns
5 train_offer_ = pd.merge(train_offer_, a[['id', 'has_bought_category_90']], how='left'
6 train_offer_.columns
```

Out[79]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bough
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_ti
mes',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60', 'has_bought_company_90', 'has_bought_b
rand_30',
      'has_bought_brand_60', 'has_bought_brand_90', 'has_bought_categ
ory_30',
      'has_bought_category_60', 'has_bought_category_90'],
      dtype='object')
```

In [80]:

```
1 print(train_offer_.chain_x.unique())
2 train_offer_.columns
```

[4 2 8]

Out[80]:

```
Index(['id', 'chain_x', 'offer', 'market', 'repeattrips', 'repeater',
      'offerdate', 'offer_category', 'quantity', 'offer_company',
      'offervalue', 'offer_brand', 'date', 'has_bought_company',
      'has_bought_company_amount', 'has_bought_company_q', 'has_bough
t_brand',
      'has_bought_brand_amount', 'has_bought_brand_q', 'has_bought_ca
tegory',
      'has_bought_category_amount', 'has_bought_category_q', 'has_bou
ght_all',
      'has_bought_all_amount', 'has_bought_all_q', 'Delta', 'total_ti
mes',
      'total_q', 'total_amount', 'has_bought_company_30',
      'has_bought_company_60', 'has_bought_company_90', 'has_bought_b
rand_30',
      'has_bought_brand_60', 'has_bought_brand_90', 'has_bought_categ
ory_30',
      'has_bought_category_60', 'has_bought_category_90'],
      dtype='object')
```

```
1 ### Finish up feature engineering
```

In [83]:

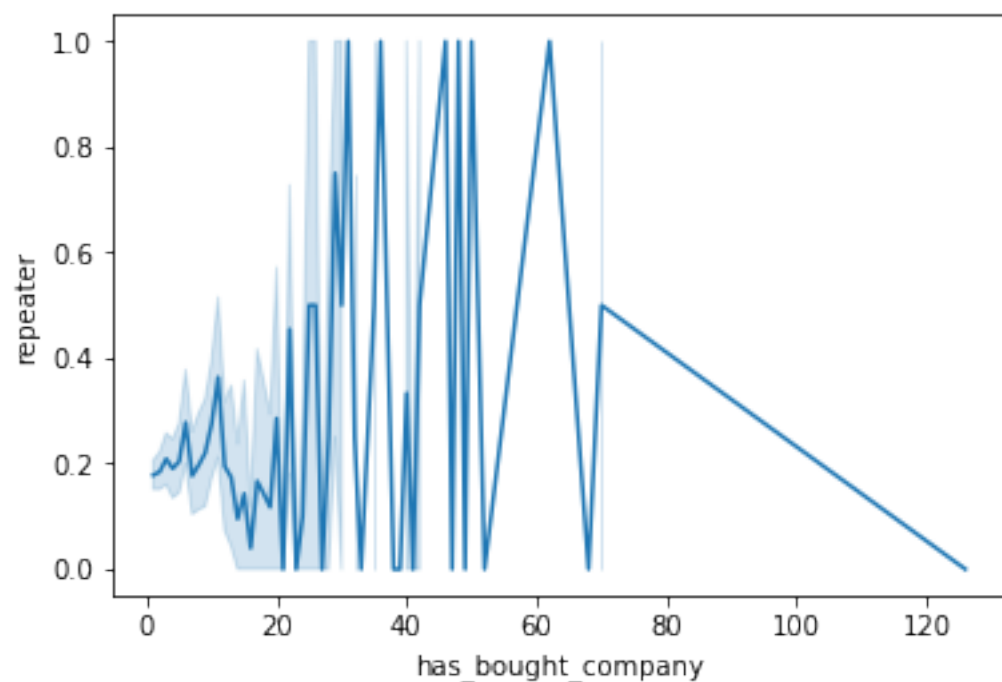
```
1 train_offer_ = train_offer_.where(train_offer_.notnull(), 0)
```

## Data exploration and modeling

Then, we use lineplot to explore the pairwise relationship between the number of companies each customer buys items from and the corresponding return rates.

In [32]:

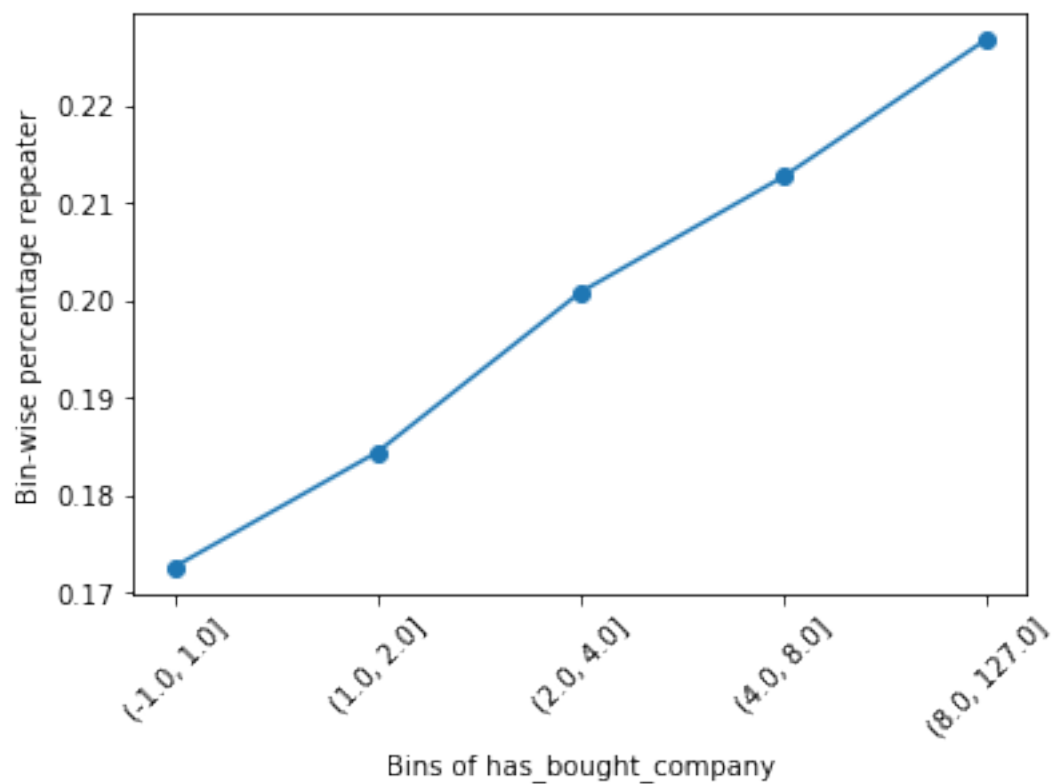
```
1 import seaborn as sns
2 ax = sns.lineplot(x = "has_bought_company",y="repeater",data = train_offer_)
```

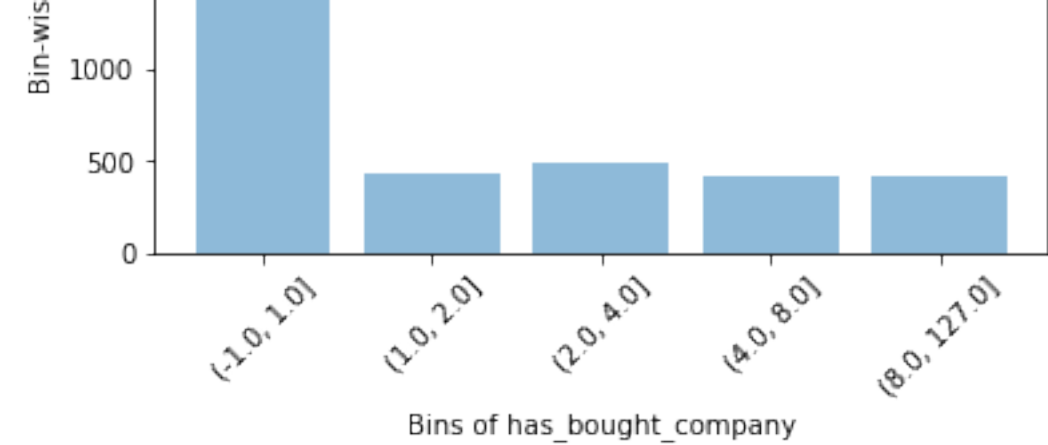


In [123]:

```
1 univariate_plotter(feature='has_bought_company', train_sub_col_target=train_offe
```

Below are the actual found rates wrt has\_bought\_company





Out[123]:

	has_bought_company	y_mean	y_sum	has_bought_company_mean
0	(-1.0, 1.0]	0.172597	3007	0.237113
1	(1.0, 2.0]	0.184397	423	2.000000
2	(2.0, 4.0]	0.200828	483	3.445135
3	(4.0, 8.0]	0.212714	409	6.210269
4	(8.0, 127.0]	0.226829	410	16.765854

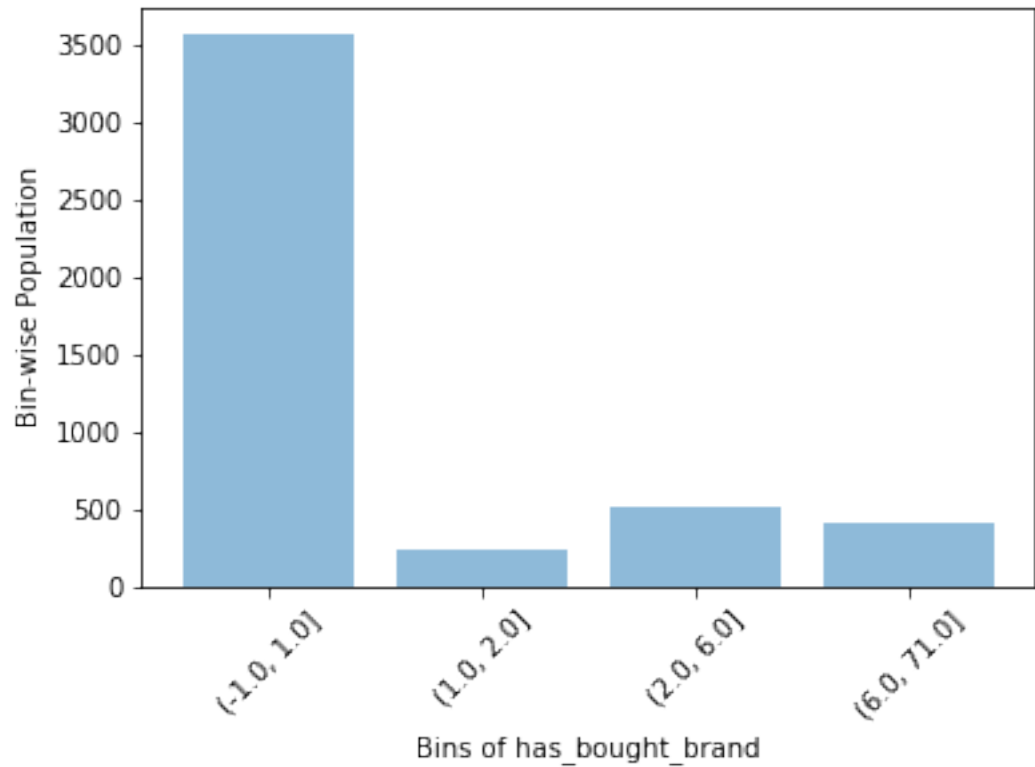
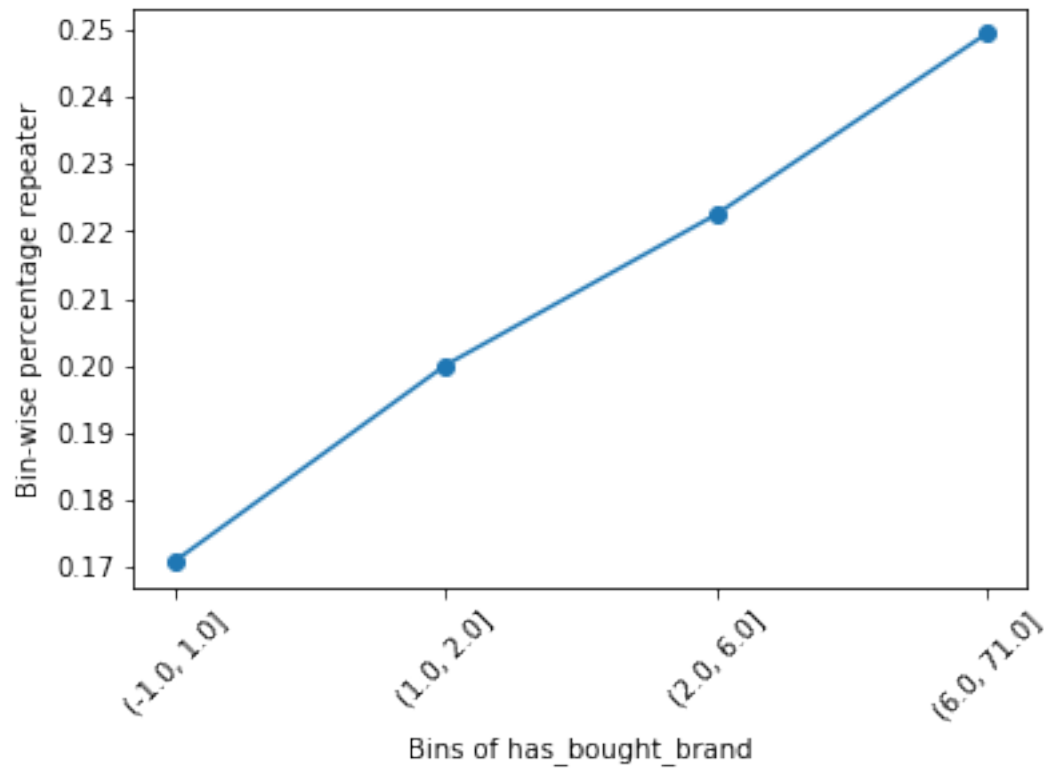
1

Next, we classify the cases into different groups by creating bins of number of brands each customer buys. This way, we can visualize the return rates and number of customers(population) by groups.

In [122]:

```
1 univariate_plotter(feature='has_bought_brand', train_sub_col_target=train_offer_
```

Below are the actual found rates wrt has\_bought\_brand



Out[122]:

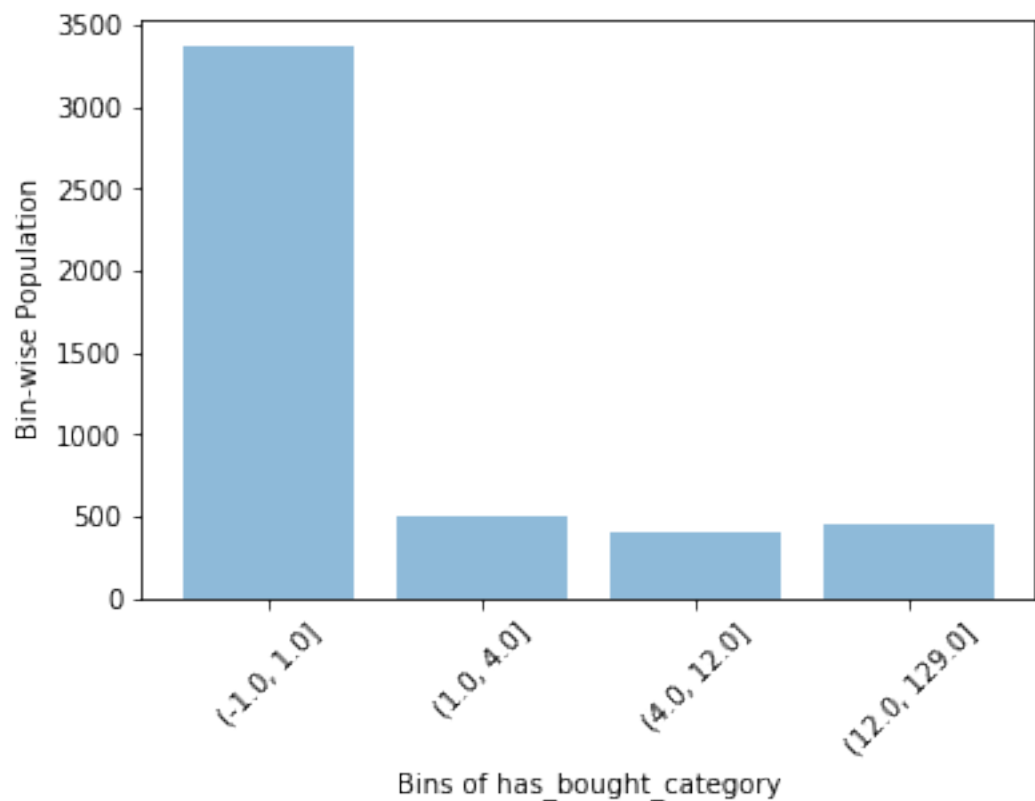
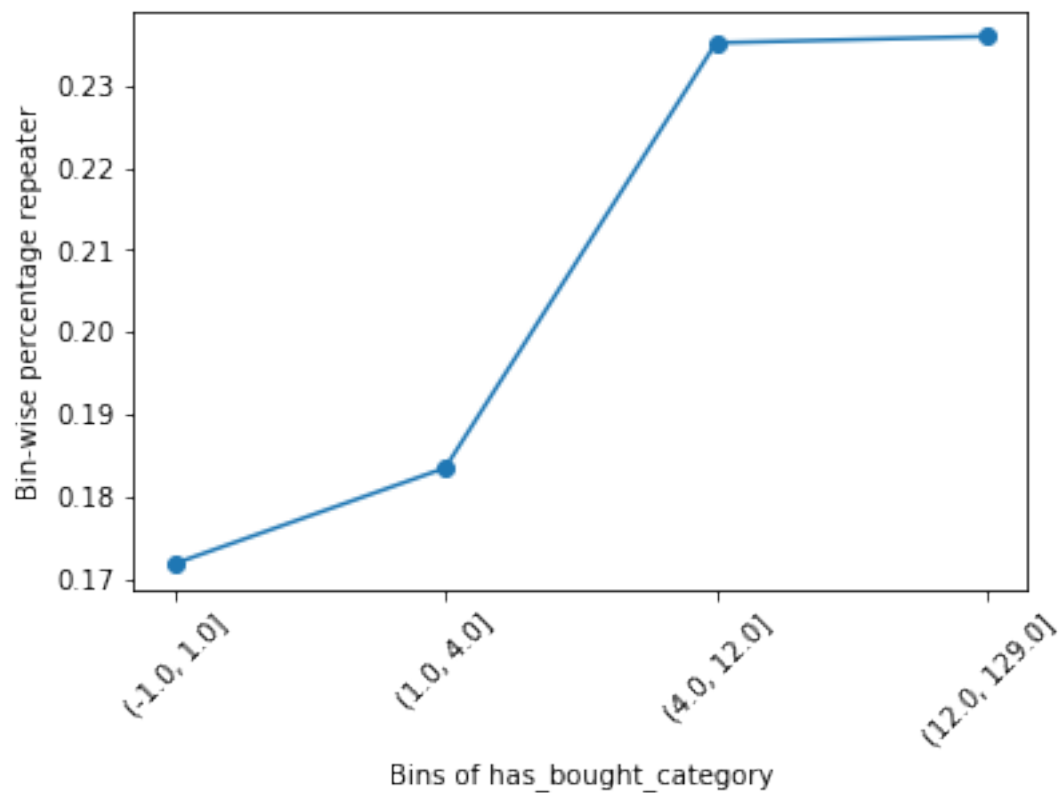
	has_bought_brand	y_mean	y_sum	has_bought_brand_mean
0	(-1.0, 1.0]	0.170827	3565	0.111360
1	(1.0, 2.0]	0.200000	245	2.000000
2	(2.0, 6.0]	0.222437	517	4.189555
3	(6.0, 71.0]	0.249383	405	14.333333



In [124]:

```
1 univariate_plotter(feature='has_bought_category', train_sub_col_target=train_of:
```

Below are the actual found rates wrt has\_bought\_category



Out[124]:

	has_bought_category	y_mean	y_sum	has_bought_category_mean
0	(-1.0, 1.0]	0.171810	3370	0.148071
1	(1.0, 4.0]	0.183468	496	2.812500
2	(4.0, 12.0]	0.235149	404	7.646040
3	(12.0, 129.0]	0.235931	462	25.614719

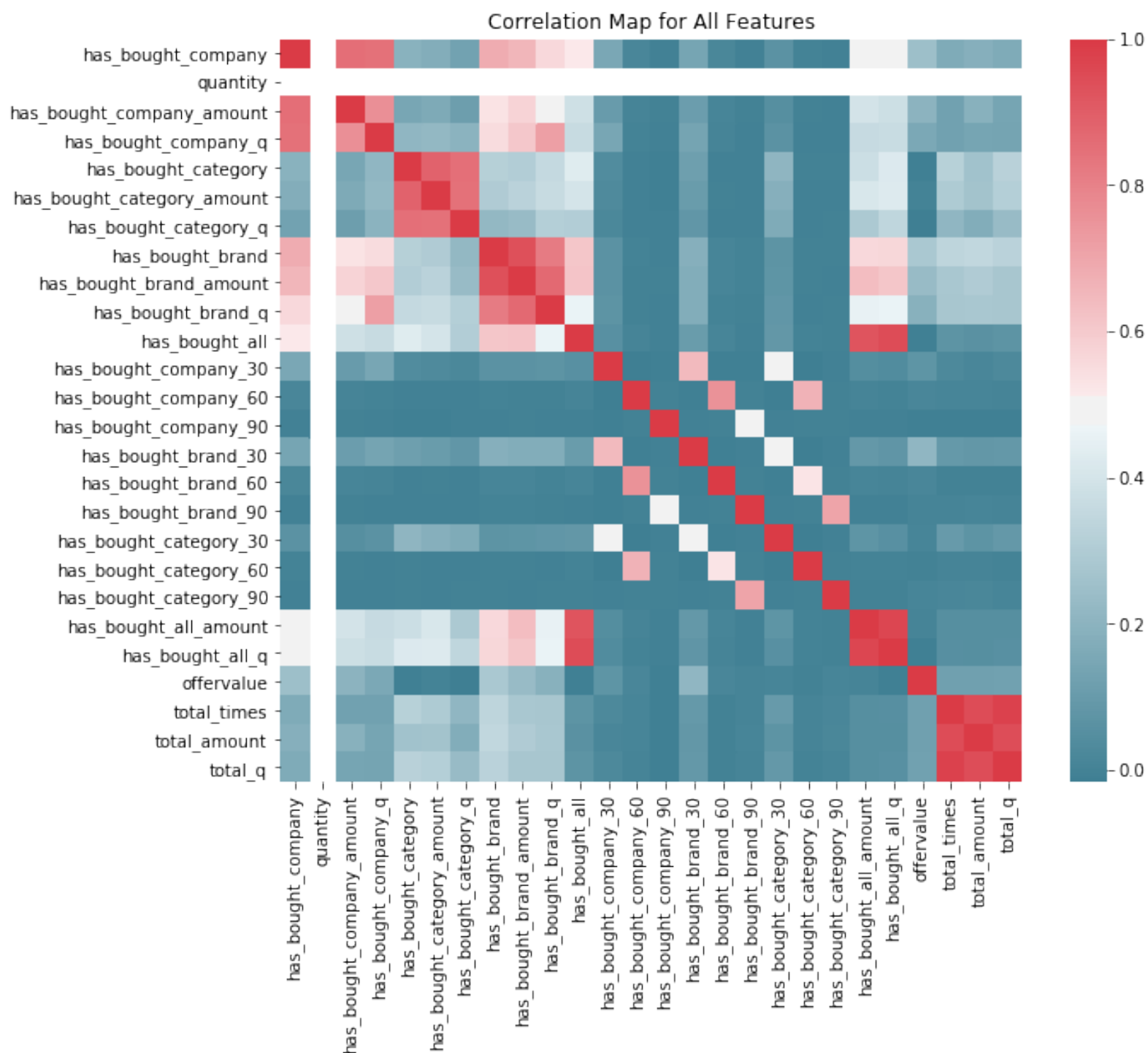
Next, we will conduct the correlation analysis among the features so that we can have a more intuitive understanding of the covariation and relationships of these features. In the Positive correlations are displayed in red and negative correlations in green color. Color intensity is proportional to the correlation coefficients.

In [84]:

[illegible]

In [119]:

```
1 import seaborn as sns
2 f, ax = plt.subplots(figsize=(10, 8))
3 corr = X_train.corr()
4 ax = sns.heatmap(corr, mask=np.zeros_like(corr, dtype=np.bool), cmap=sns.diverg
5                 square=True, ax=ax)
6 ax.set_title('Correlation Map for All Features')
7 plt.show()
```



## Model Building

### 1. Logistic Regression

[illegible]

```
1 from sklearn.grid_search import GridSearchCV
2 from sklearn import grid_search
3 from sklearn.metrics import roc_auc_score
4 logreg = LogisticRegression(penalty='l1', solver='liblinear', C=100)
5 logreg.fit(X_train, y_train)
6 y_pred = logreg.predict_proba(X_test)[:,-1]
7 print('Accuracy of logistic regression classifier on training set: {:.2f}'.format(
8 print('Accuracy of logistic regression classifier on test set: {:.2f}'.format(
9 performance = roc_auc_score(y_test, y_pred)
10 print('Logistic Regression: Area under the ROC curve = {}'.format(performance))
```

```
1 logreg = LogisticRegression()
2 rfe = RFE(logreg, 18)
3 rfe = rfe.fit(X_train, y_train)
4 print(rfe.support_)
5 print(rfe.ranking_)

[ True  True False  True  True False False False  True  True  True  Tr
ue
  True  True  True  True  True  True  True  True False  True  True Fal
se
  False False]
[1 1 6 1 1 3 4 2 1 1 1 1 1 1 1 1 1 1 5 1 1 7 9 8]
```

In [59]:

```
1 from sklearn.feature_selection import f_classif
2 f_classif(X_train,y_train)
```

C:\Users\kathl\Anaconda3\lib\site-packages\sklearn\feature\_selection\univariate\_selection.py:113: UserWarning: Features [1] are constant.

UserWarning)

C:\Users\kathl\Anaconda3\lib\site-packages\sklearn\feature\_selection\univariate\_selection.py:114: RuntimeWarning: invalid value encountered

in true\_divide

f = msb / msb

Out[59]:

```
(array([ 4.04023787,          nan,  3.15344253,  1.46101095,
        10.96661338,  19.69761165,  10.70177851,  10.17324837,
        13.48587595,  4.66126931,  29.99636472,  0.20566592,
         0.07314886,  0.44593301,  0.86035822,  0.12306138,
         0.44593301,  2.99475209,  0.17648542,  0.22289444,
        31.48799709,  29.73017011,  0.50925039,  10.05525484,
        13.72673972,  15.70253885]),
 array([ 4.44984238e-02,          nan,  7.58473297e-02,
         2.26844880e-01,  9.36389104e-04,  9.33001020e-06,
         1.07992686e-03,  1.43656805e-03,  2.43688604e-04,
         3.09130577e-02,  4.61025881e-08,  6.50211842e-01,
         7.86820782e-01,  5.04313392e-01,  3.53697639e-01,
         7.25757137e-01,  5.04313392e-01,  8.36162448e-02,
         6.74435031e-01,  6.36870525e-01,  2.15067051e-08,
         5.28304667e-08,  4.75507557e-01,  1.53128306e-03,
         2.14444721e-04,  7.54918789e-05]))
```

In [60]:

```
1 import statsmodels.api as sm
2 logit_model=sm.Logit(y_train.astype(float),X_train.astype(float))
3 result=logit_model.fit()
4 print(result.summary())
```

Warning: Maximum number of iterations has been exceeded.

Current function value: 0.466863

Iterations: 35

#### Logit Regression Results

=====

Dep. Variable: repeater No. Observations:

3785

Model: Logit Df Residuals:

3759

Method: MLE Df Model:

25

Date: Fri, 29 Jun 2018 Pseudo R-squ.:

0.01684

Time: 09:41:08 Log-Likelihood:

-1767.1				
converged:	False	LL-Null:		
-1797.3				
		LLR p-value:		8
.787e-05				
=====				
=====				
	coef	std err	z	P> z
[95.0% Conf. Int.]				
-----				
-----				
has_bought_company	-0.0157	0.033	-0.477	0.633
-0.080 0.049				
quantity	-1.6267	0.118	-13.732	0.000
-1.859 -1.394				
has_bought_company_amount	-0.0009	0.003	-0.312	0.755
-0.006 0.005				
has_bought_company_q	0.0076	0.016	0.492	0.623
-0.023 0.038				
has_bought_category	-0.0167	0.012	-1.387	0.166
-0.040 0.007				
has_bought_category_amount	0.0031	0.002	1.357	0.175
-0.001 0.008				
has_bought_category_q	0.0030	0.005	0.637	0.524
-0.006 0.012				
has_bought_brand	-0.0100	0.038	-0.266	0.790
-0.084 0.064				
has_bought_brand_amount	0.0096	0.007	1.373	0.170
-0.004 0.023				
has_bought_brand_q	-0.0287	0.022	-1.329	0.184
-0.071 0.014				
has_bought_all	0.0940	0.071	1.332	0.183
-0.044 0.232				
has_bought_company_30	-0.0439	0.194	-0.226	0.821
-0.424 0.336				
has_bought_company_60	-11.3514	426.438	-0.027	0.979
-847.155 824.452				
has_bought_company_90	-23.7044	3.1e+05	-7.65e-05	1.000
-6.07e+05 6.07e+05				
has_bought_brand_30	0.0211	0.249	0.085	0.932
-0.466 0.508				
has_bought_brand_60	12.2168	426.439	0.029	0.977
-823.589 848.023				
has_bought_brand_90	-7.0763	1.03e+07	-6.87e-07	1.000
-2.02e+07 2.02e+07				
has_bought_category_30	0.2066	0.189	1.094	0.274
-0.164 0.577				
has_bought_category_60	-0.5952	1.511	-0.394	0.694
-3.557 2.366				
has_bought_category_90	-12.7568	1.03e+07	-1.24e-06	1.000
-2.02e+07 2.02e+07				
has_bought_all_amount	-0.0002	0.014	-0.012	0.991
-0.028 0.028				

has_bought_all_q	-0.0081	0.058	-0.140	0.889
-0.121	0.105			
offervalue	-0.0494	0.092	-0.535	0.593
-0.231	0.132			
total_times	-0.0006	0.000	-2.781	0.005
-0.001	-0.000			
total_amount	3.348e-05	3.16e-05	1.058	0.290
-2.85e-05	9.55e-05			
total_q	0.0004	0.000	3.004	0.003
0.000	0.001			

=====

=====

```
C:\Users\kathl\Anaconda3\lib\site-packages\statsmodels\base\model.py:466: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle_retvals
  "Check mle_retvals", ConvergenceWarning)
```

## 2.Decision Tree

In [96]:

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn import tree
3 clf = tree.DecisionTreeClassifier(criterion = "gini", splitter = 'random', max_
4 clf = clf.fit(X_train, y_train)
5 y_pred_tree = clf.predict_proba(X_test)[: ,1]
6 print('Accuracy of Decision Tree classifier on training set: {:.2f}'.format(clf
7 print('Accuracy of Decision Tree classifier on test set: {:.2f}'.format(clf.sco
8 print(clf.feature_importances_)
9 performance = roc_auc_score(y_test, y_pred_tree)
10 print ('DecisionTree: Area under the ROC curve = {}'.format(performance))
```

```
Accuracy of Decision Tree classifier on training set: 0.82
Accuracy of Decision Tree classifier on test set: 0.81
[0.          0.          0.          0.          0.          0.
 0.08583677 0.          0.          0.          0.119805 0.
 0.          0.          0.          0.          0.          0.
 0.          0.          0.          0.69656351 0.          0.
 0.          0.09779472]
DecisionTree: Area under the ROC curve = 0.593498917317226
```

In [80]:

```
1 rfc = RandomForestClassifier(n_jobs=-1, max_features='sqrt', oob_score = True)
2 # Use a grid over parameters of interest
3 param_grid = {
4     "n_estimators" : [9, 18, 27, 36, 45, 54, 63],
5     "max_depth" : [1, 5, 10, 15, 20, 25, 30],
6     "min_samples_leaf" : [1, 2, 4, 6, 8, 10]}
7
8 CV_rfc = GridSearchCV(estimator=rfc, param_grid=param_grid, cv= 10)
9 CV_rfc.fit(X_train, y_train)
10 print(CV_rfc.best_params_)
```

C:\Users\kathl\Anaconda3\lib\site-packages\sklearn\ensemble\forest.py:439: UserWarning: Some inputs do not have OOB scores. This probably means too few trees were used to compute any reliable oob estimates.

warn("Some inputs do not have OOB scores. ")

C:\Users\kathl\Anaconda3\lib\site-packages\sklearn\ensemble\forest.py:444: RuntimeWarning: invalid value encountered in true\_divide

predictions[k].sum(axis=1)[: , np.newaxis])

C:\Users\kathl\Anaconda3\lib\site-packages\sklearn\ensemble\forest.py:439: UserWarning: Some inputs do not have OOB scores. This probably means too few trees were used to compute any reliable oob estimates.

warn("Some inputs do not have OOB scores. ")

C:\Users\kathl\Anaconda3\lib\site-packages\sklearn\ensemble\forest.py:444: RuntimeWarning: invalid value encountered in true\_divide

predictions[k].sum(axis=1)[: , np.newaxis])

C:\Users\kathl\Anaconda3\lib\site-packages\sklearn\ensemble\forest.py:439: UserWarning: Some inputs do not have OOB scores. This probably means too few trees were used to compute any reliable oob estimates.

warn("Some inputs do not have OOB scores. ")

C:\Users\kathl\Anaconda3\lib\site-packages\sklearn\ensemble\forest.py:444: RuntimeWarning: invalid value encountered in true\_divide

### 3. Random Forest



In [112]:

```
1  from sklearn.ensemble import RandomForestClassifier
2  # Instantiate model with 1000 decision trees
3  rf = RandomForestClassifier(criterion = "gini", max_leaf_nodes= 18,min_samples_
4  # min_samples_leaf = 10, 15
5  # Train the model on training data
6  rf.fit(X_train, y_train);
7
8  predictions = rf.predict_proba(X_test)[: ,1]
9  print('Accuracy of Random Forest classifier on training set: {:.2f}'.format(rf.s
10 print('Accuracy of Random Forest classifier on test set: {:.2f}'.format(rf.score
11
12 performance = roc_auc_score(y_test, predictions)
13 print ('Random Forest: Area under the ROC curve = {}'.format(performance))
```

Accuracy of Random Forest classifier on training set: 0.83

Accuracy of Random Forest classifier on test set: 0.81

Random Forest: Area under the ROC curve = 0.6122001253632686

In [100]:

```
1  print(rf.feature_importances_)
2  print(cols)
```

```
[ 0.04201335  0.          0.05516629  0.04264195  0.0412595   0.036087
02
 0.03503255  0.02968182  0.02832364  0.0273683   0.08371838  0.068406
45
 0.11483007  0.07889876  0.0210911   0.07277819  0.10214323  0.120559
4 ]
['has_bought_company', 'quantity', 'has_bought_company_amount', 'has_b
ought_company_q', 'has_bought_category', 'has_bought_category_amount',
'has_bought_category_q', 'has_bought_brand', 'has_bought_brand_amount'
, 'has_bought_brand_q', 'Delta', 'has_bought_all', 'has_bought_all_amo
unt', 'has_bought_all_q', 'offervalue', 'total_times', 'total_amount',
'total_q']
```

## 4. XGboost

In [100]:

```
1 import xgboost as xgb
2 dtrain=xgb.DMatrix(X_train,label=y_train)
3 dtest=xgb.DMatrix(X_test)
4
5 params={'booster':'gbtree',
6         'objective': 'binary:logistic',
7         'eval_metric': 'auc',
8         'max_depth':3,
9         'gamma':0,
10        'lambda':3,
11        'subsample':0.8,
12        'colsample_bytree':0.8,
13        'min_child_weight':3,
14        'eta': 0.007,
15        'seed':0,
16        'nthread':8,
17        'silent':0}
18
19 watchlist = [(dtrain, 'train')]
```

In [103]:

```
1 bst=xgb.train(params,dtrain,num_boost_round=500,evals=watchlist)
2 ypred=bst.predict(dtest)
3 y_pred = (ypred >= 0.5)*1
```

```
[22:55:13] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 12
extra nodes, 0 pruned nodes, max_depth=3
[0]      train-auc:0.578392
[22:55:13] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 14
extra nodes, 0 pruned nodes, max_depth=3
[1]      train-auc:0.587779
[22:55:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 12
extra nodes, 0 pruned nodes, max_depth=3
[2]      train-auc:0.591104
[22:55:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 10
extra nodes, 0 pruned nodes, max_depth=3
[3]      train-auc:0.592512
[22:55:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 12
extra nodes, 0 pruned nodes, max_depth=3
[4]      train-auc:0.592114
[22:55:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 10
extra nodes, 0 pruned nodes, max_depth=3
[5]      train-auc:0.592253
[22:55:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 14
extra nodes, 0 pruned nodes, max_depth=3
```

In [115]:

```
1 from sklearn import metrics
2 print('XGBoost: Area under the ROC curve = %.4f' % metrics.roc_auc_score(y_test,
```

XGBoost: Area under the ROC curve = 0.6288

In [110]:

```
1 from xgboost import XGBClassifier
2 model = XGBClassifier()
3 model.fit(X_train, y_train)
4 # feature importance
5 print(model.feature_importances_)
6 importances = list(model.feature_importances_)
7 feature_list = list(cols)
8 feature_importances = [(feature,round(importance,2)) for feature,importance in zip(feature_list,importances)]
9 feature_importances = sorted(feature_importances,key=lambda x:x[1],reverse = True)
10 [print('Variable: {:20}Importance: {}'.format(*pair))for pair in feature_importances]
```

```
[0.0374415  0.          0.10140406 0.05928237 0.03432137 0.07176287
 0.02964118 0.01716069 0.03120125 0.01560062 0.00468019 0.00468019
 0.          0.          0.00156006 0.          0.          0.00624025
 0.          0.          0.05460218 0.00156006 0.03276131 0.12948518
 0.18876755 0.17784712]
```

```
Variable: total_amount      Importance: 0.1899999976158142
Variable: total_q           Importance: 0.18000000715255737
Variable: total_times       Importance: 0.12999999523162842
Variable: has_bought_company_amountImportance: 0.10000000149011612
Variable: has_bought_category_amountImportance: 0.07000000029802322
Variable: has_bought_company_qImportance: 0.05999999865889549
Variable: has_bought_all_amountImportance: 0.05000000074505806
Variable: has_bought_company Importance: 0.03999999910593033
Variable: has_bought_category Importance: 0.029999999329447746
Variable: has_bought_category_qImportance: 0.029999999329447746
Variable: has_bought_brand_amountImportance: 0.029999999329447746
Variable: offervalue        Importance: 0.029999999329447746
Variable: has_bought_brand   Importance: 0.019999999552965164
Variable: has_bought_brand_q Importance: 0.019999999552965164
Variable: has_bought_category_30Importance: 0.009999999776482582
Variable: quantity          Importance: 0.0
Variable: has_bought_all     Importance: 0.0
Variable: has_bought_company_30Importance: 0.0
Variable: has_bought_company_60Importance: 0.0
Variable: has_bought_company_90Importance: 0.0
Variable: has_bought_brand_30 Importance: 0.0
Variable: has_bought_brand_60 Importance: 0.0
Variable: has_bought_brand_90 Importance: 0.0
Variable: has_bought_category_60Importance: 0.0
Variable: has_bought_category_90Importance: 0.0
Variable: has_bought_all_q   Importance: 0.0
```

