

STARBUCKS CAPSTONE

Sardor Israilov

UDACITY

Content

Introduction and Domain background	2
Problem statement	2
Solution:	2
Starbucks Ipython:	2
Solution statement	2
Evaluation metrics	3
Questions:	3
Results for binary classification on the validation set	3
FINDINGS	3
Income	3
Gender Analysis	4
Time analysis:	5
Transaction:	5
Responses:	5
Project design	5
Annexes	6
Story of 1 person:	6
Datasets and inputs	6
References:	8

Introduction and Domain background

The project was an experiment from Starbucks to optimize the offers sent to its customers and to determine what is the best offer on the individual level. Some people may respond negatively to the offer, so in some cases it's best not to send any offer at all.

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). The project is meant to optimize the offers sent to the customers based on the certain characteristics as sex, income range, demographics etc.

Example

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.

Problem statement

The problem that I intend to solve consists of determining:

whether the customer will view and complete the offer or just view it.

The most crucial step would be to understand the data and its probable implications on the outcome. Data pre-processing would be the crucial step to the further results. After obtaining good results, the next step would be to identify the crucial factors for the offer to be completed and indicate the strategies for the marketing campaign: which offer to send to which person.

Solution:

The project is about data analysis and machine learning. At first, we respond to the questions that would be interesting from business point of view, afterwards we construct machine learning model for prediction of offer quality.

Starbucks jupyter Ipython:

Steps taken during the project are:

- data exploration and visualization **Data_Analysis.ipynb**
- preprocessing and cleaning
- merging in 1 dataset
- preparing train,test data and labels
- different classifiers and their results
- Bayesian optimization for the best- RFC **ML_SUPERVISED.ipynb**

Solution statement

The first thing to do would be to understand and visualize the data. Then, we would eliminate the outliers, the (NaN) rows in the table that would not contribute to our training or test. Afterwards, I would combine 3 data tables in a big one with preprocessing steps. For example, instead of having male, female

we would have 0 or 1. Use data scaling as the way to normalize the input data. Afterwards, we would try several models of Sklearn and Pytorch CNN and choose the best one.

Evaluation metrics

The metric that was primarily used for binary classification is Accuracy. $\frac{TP+TN}{TP+TN+FP+FN}$. Although, in the test sample we have 8589 and 1709 positives, with large number of negative samples precision is better. [\[1\]](#)

Questions TO RESPOND :

- 1- What are the customers that don't respond to offers i.e never view them?
- 2- Based on the demographic data of the customers who gets the highest income range , males or females?
- 3- What are min,max,average transaction.
- 4- Which year/month got most clients for starbucks?
- 5- Which type of promotions(offers) each gender likes?
- 6- What is the average length between two transcript for the same customer?

Results for binary classification on the validation set

	XGBoost [3]	RFC [2]	LinearSVC [7]	KNearestN	DecisionTreeClassifier
accuracy	0.81	0.82	0.67	0.52	0.77
precision	0.79	0.81	0.47		
recall	0.84	0.84	0.78		

We can predict with 81-84% accuracy if the offer will be completed in the future. Bayesian optimization[\[5\]](#) was performed on the RFC.

FINDINGS

Income

On the income range people equally(33% for each) fall into 3 categories:

Low income 0-55000

Average: 55000-75000

ISRAILOV SARDOR

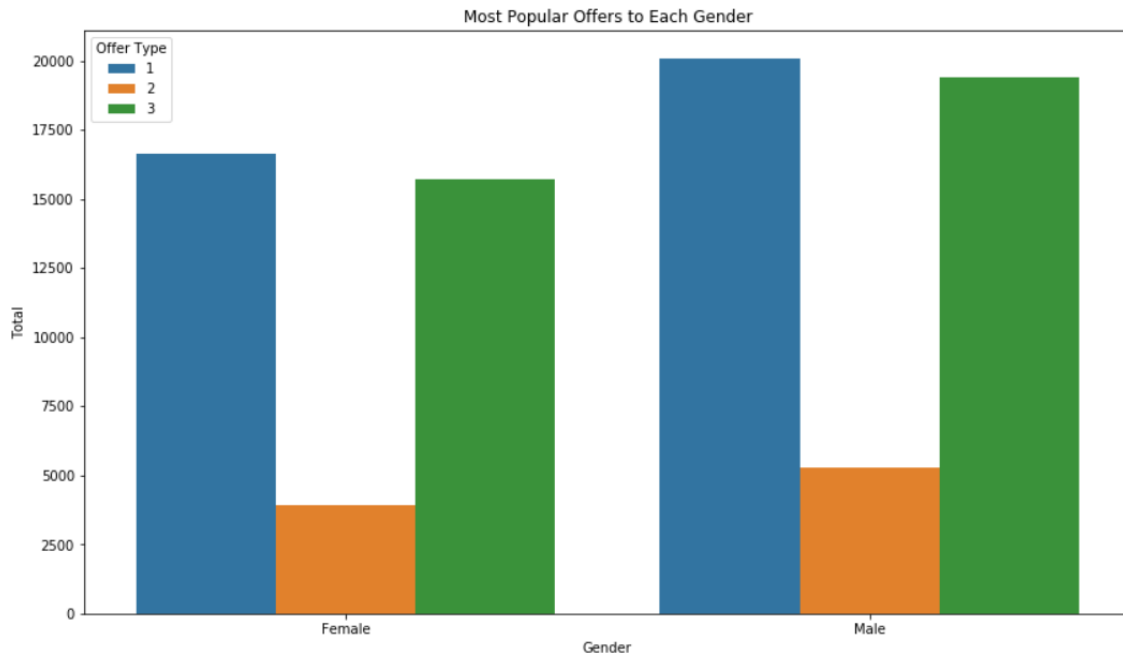
High: 75000-120000

The average is 64.000

Gender Analysis

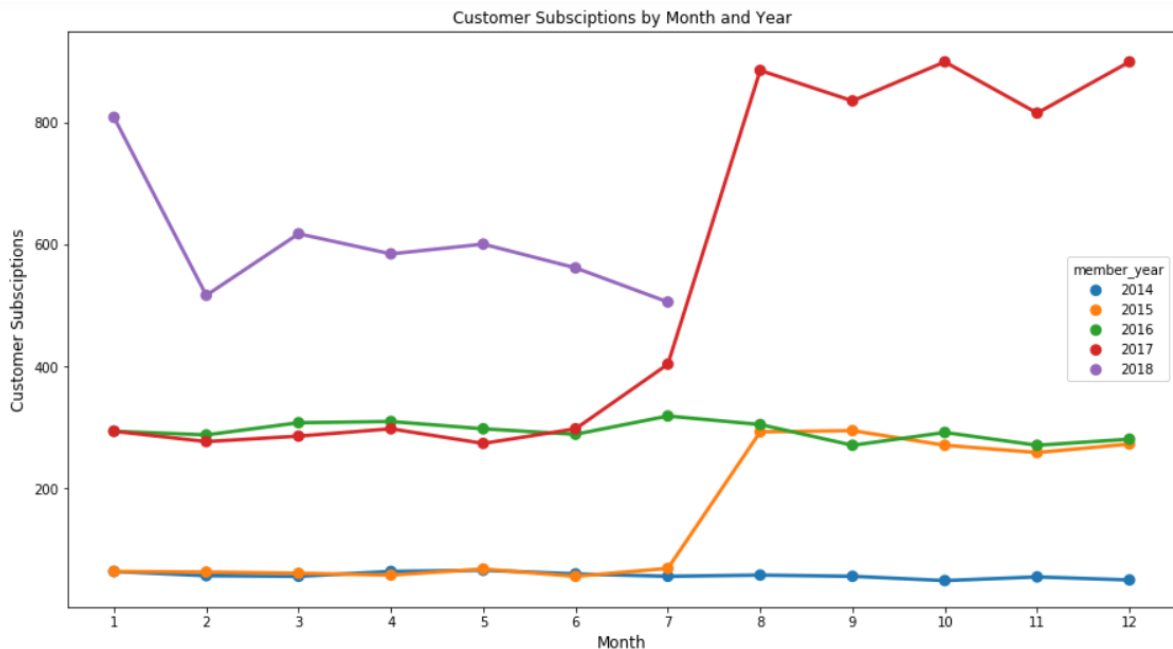
57.2% of starbucks customers are males with 4.3% as female and 1.43 as Other.

In, general high-income female visits most starbucks.



We see that BOGO and Discount are most popular for both male and female.

Time analysis:



Most of the customers came during 2017 Year. Year 217 and 2016 suggest that some interesting trend is happening during august to December (customer increase).

Transaction:

12.8-mean

0.05 -min

1062-max

Responses:

- 1) Number of persons that don't open an offer 2255
- 2) Females

Project design

I have identified the following steps to project resolution:

- Visualising&Data analysis
- Pre-processing data
- Scaling the numerical features
- Testing supervised models
- Evaluating the results on the accuracy. If the results to be improves use GridSearchSV or Bayesian hyper-parameter tuning.
- Conclusion and evaluation of the crucial factors for the decision i.e. which columns contribute most for the decision of whether the customer will view&complete an offer or just view it.

Conclusion:

It was demonstrated hereby, that 2175 customers never open their offers, so it might be a good idea to put them off list. Gender, year and transaction analysis were discussed. The basic classifier was constructed that could predict with 84% accuracy whether the offer would be completed or not (before cross-validation). I would personally adapt Bayesian search for XGBoost and choose it for this task.

Annexe:

Story of 1 person:

```
In [113]: transcript[transcript.person=='78afa995795e4d85b5d9ceeca43f5fef']
```

Out[113]:		event	person	time	value
	0	offer received	78afa995795e4d85b5d9ceeca43f5fef	0	{'offer id': '9b98b8c7a33c4b65b9aebfe6a799e6d9'}
	15561	offer viewed	78afa995795e4d85b5d9ceeca43f5fef	6	{'offer id': '9b98b8c7a33c4b65b9aebfe6a799e6d9'}
	47582	transaction	78afa995795e4d85b5d9ceeca43f5fef	132	{'amount': 19.89}
	47583	offer completed	78afa995795e4d85b5d9ceeca43f5fef	132	{'offer_id': '9b98b8c7a33c4b65b9aebfe6a799e6d9'}
	49502	transaction	78afa995795e4d85b5d9ceeca43f5fef	144	{'amount': 17.78}
	53176	offer received	78afa995795e4d85b5d9ceeca43f5fef	168	{'offer id': '5a8bc65990b245e5a138643cd4eb9837'}
	85291	offer viewed	78afa995795e4d85b5d9ceeca43f5fef	216	{'offer id': '5a8bc65990b245e5a138643cd4eb9837'}
	87134	transaction	78afa995795e4d85b5d9ceeca43f5fef	222	{'amount': 19.67}
	92104	transaction	78afa995795e4d85b5d9ceeca43f5fef	240	{'amount': 29.72}
	141566	transaction	78afa995795e4d85b5d9ceeca43f5fef	378	{'amount': 23.93}
	150598	offer received	78afa995795e4d85b5d9ceeca43f5fef	408	{'offer id': 'ae264e3637204a6fb9bb56bc8210ddfd'}
	163375	offer viewed	78afa995795e4d85b5d9ceeca43f5fef	408	{'offer id': 'ae264e3637204a6fb9bb56bc8210ddfd'}
	201572	offer received	78afa995795e4d85b5d9ceeca43f5fef	504	{'offer id': 'f19421c1d4aa40978ebb69ca19b0e20d'}
	218393	transaction	78afa995795e4d85b5d9ceeca43f5fef	510	{'amount': 21.72}
	218394	offer completed	78afa995795e4d85b5d9ceeca43f5fef	510	{'offer_id': 'ae264e3637204a6fb9bb56bc8210ddfd'}
	218395	offer completed	78afa995795e4d85b5d9ceeca43f5fef	510	{'offer_id': 'f19421c1d4aa40978ebb69ca19b0e20d'}
	230412	transaction	78afa995795e4d85b5d9ceeca43f5fef	534	{'amount': 26.56}
	262138	offer viewed	78afa995795e4d85b5d9ceeca43f5fef	582	{'offer id': 'f19421c1d4aa40978ebb69ca19b0e20d'}

Datasets and inputs

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) - type of offer ie BOGO, discount, informational
- difficulty (int) - minimum required spend to complete an offer

ISRAILOV SARDOR

- reward (int) - reward given for completing an offer
- duration (int) - time for offer to be open, in days
- channels (list of strings)

profile.json

- age (int) - age of the customer
- became_member_on (int) - date when customer created an app account
- gender (str) - gender of the customer (note some entries contain 'O' for other rather than M or F)
- id (str) - customer id
- income (float) - customer's income

transcript.json

- event (str) - record description (ie transaction, offer received, offer viewed, etc.)
- person (str) - customer id
- time (int) - time in hours since start of test. The data begins at time $t=0$
- value - (dict of strings) - either an offer id or transaction amount depending on the record

In [2]: `portfolio.head()`

Out[2]:

	channels	difficulty	duration	id	offer_type	reward
0	[email, mobile, social]	10	7	ae264e3637204a6fb9bb56bc8210ddfd	bogo	10
1	[web, email, mobile, social]	10	5	4d5c57ea9a6940dd891ad53e9dbe8da0	bogo	10
2	[web, email, mobile]	0	4	3f207df678b143eea3cee63160fa8bed	informational	0
3	[web, email, mobile]	5	7	9b98b8c7a33c4b65b9aebfe6a799e6d9	bogo	5
4	[web, email]	20	10	0b1e1539f2cc45b7b9fa7c272da2e1d7	discount	5

In [3]: `profile.head()`

Out[3]:

	age	became_member_on	gender	id	income
0	118	20170212	None	68be06ca386d4c31939f3a4f0e3dd783	NaN
1	55	20170715	F	0610b486422d4921ae7d2bf64640c50b	112000.0
2	118	20180712	None	38fe809add3b4fcf9315a9694bb96ff5	NaN
3	75	20170509	F	78afa995795e4d85b5d9ceeca43f5fef	100000.0
4	118	20170804	None	a03223e636434f42ac4c3df47e8bac43	NaN

In [4]: `transcript.head()`

Out[4]:

	event	person	time	value
0	offer received	78afa995795e4d85b5d9ceeca43f5fef	0	{'offer id': '9b98b8c7a33c4b65b9aebfe6a799e6d9'}
1	offer received	a03223e636434f42ac4c3df47e8bac43	0	{'offer id': '0b1e1539f2cc45b7b9fa7c272da2e1d7'}
2	offer received	e2127556f4f64592b11af22de27a7932	0	{'offer id': '2906b810c7d4411798c6938adc9daaa5'}
3	offer received	8ec6ce2a7e7949b1bf142def7d0e0586	0	{'offer id': 'fafcd668e3743c1bb461111dcafc2a4'}
4	offer received	68617ca6246f4fbc85e91a2a49552598	0	{'offer id': '4d5c57ea9a6940dd891ad53e9dbe8da0'}

References:

- 1) <https://towardsdatascience.com/what-metrics-should-we-use-on-imbalanced-data-set-precision-recall-roc-e2e79252aeba>
- 2) <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
- 3) <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingClassifier.htmlhttps://www.kdnuggets.com/2019/07/xgboost-random-forest-bayesian-optimisation.html>
- 4) <https://mlfromscratch.com/gridsearch-keras-sklearn/#/>
- 5) <https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html>
- 6) <https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.merge.html>
- 7) <https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html>