

Machine Learning Engineering Nanodegree

Capstone Project

Mahlong Makwele Wishbert
03 March 2022

Starbucks

Definition

Project Overview

Being able to tell how customers will react to products or offers is an important part of business, StarBucks has set out to figure out just how to do this. In this project I analyze simulated StarBucks data provided by Udacity to try and figure out how customers would respond to the different offers given by StarBucks.

Problem Statement

The aim is to determine how certain customers will respond to certain offers. The data can be tricky because you have to keep in account that some customers can make purchases through the app without having received or seen the offer and thus these customers are not influenced by the offer. From a profit perspective it is also good to look at what customers will purchase still, without any offers available. You can assume that when a customer sees an offer the offer has an effect on the customer for the duration of its validation time and from this we can also see how customers would react due to the offer.

How I will deal with the problem is as follows:

- I will create a jupyter-notebook workspace

- I will import and try to merge to data files
- I will do an exploratory data analyses to try and find useful and meaningful insights
- I will train the benchmark model and other models
- I will evaluate the models

Metrics

All the models used in this problem classify data with regards to offer completion. A simple way to evaluate the models is by calculating for precision, recall and F1 score.

$$\text{Precision}(class = a) = \frac{TP(class = a)}{TP(class = a) + FP(class = a)}$$

$$\text{Recall}(class = a) = \frac{TP(class = a)}{TP(class = a) + FN(class = a)} =$$

$$\text{F-1 Score}(class = a) = \frac{2 \times \text{Precision}(class = a) \times \text{Recall}(class = a)}{\text{Precision}(class = a) + \text{Recall}(class = a)} :$$

The above equations are for calculating the metrics by which we will evaluate. The python code does the calculations for us. The data used in this project is imbalanced and therefore accuracy will not be a good enough measure since it does not distinguish between the number of correctly classified points from their respective classes. I will use the f1 score as a measure of how good the model is, because f1 will penalize the class with the most occurrence.

Analysis

Data Exploration

The Starbucks dataset comes in three files, files that need to be joined together to be able to gain insights from. The files are portfolio.json, profile.json and transcript.json. The portfolio.json file contains all the different offers that can be made and details about the offers, the data columns are as follows:

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

	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

Figure 1

The profile.json file contains information about the customer demographic, the data columns are as follows:

- 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

	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

Figure 2

The transaction.json file contains information about events when offers were received and completed, the times of completion. The transaction data columns are as follows:

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

	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': '2906b810c7d4411798c6938adc9daa5'}
3	offer received	8ec6ce2a7e7949b1bf142def7d0e0586	0	{'offer id': 'fafdcd668e3743c1bb461111dcafc2a4'}
4	offer received	68617ca6246f4fbc85e91a2a49552598	0	{'offer id': '4d5c57ea9a6940dd891ad53e9dbe8da0'}

Figure 3

I merged and combined the data above, after I had cleaned and prepared the data for machine learning algorithms I decided to take a look at the distribution of the data. I used the pandas describe method to get the data's description and it is shown below in figure 4. Looking at figure 4 I will decide to remove the email column from the data, it has a standard deviation of 0 thus it does not vary at all. This point is further justified by looking at figure 5 which shows that the correlation of the email column with the target column is 'NaN'

	age	income	reward	difficulty	duration	gender_n	web	email	mobile	social
count	44181.000000	44181.000000	44181.000000	44181.000000	44181.000000	44181.000000	44181.000000	44181.0	44181.000000	44181.000000
mean	54.301193	65380.095516	5.255517	9.647450	7.259931	0.601616	0.874064	1.0	0.873520	0.624047
std	17.378322	21643.876797	2.990987	4.456791	1.786645	0.516830	0.331782	0.0	0.332393	0.484374
min	18.000000	30000.000000	2.000000	5.000000	5.000000	0.000000	0.000000	1.0	0.000000	0.000000
25%	42.000000	49000.000000	3.000000	7.000000	7.000000	0.000000	1.000000	1.0	1.000000	0.000000
50%	55.000000	64000.000000	5.000000	10.000000	7.000000	1.000000	1.000000	1.0	1.000000	1.000000
75%	66.000000	80000.000000	10.000000	10.000000	10.000000	1.000000	1.000000	1.0	1.000000	1.000000
max	101.000000	120000.000000	10.000000	20.000000	10.000000	2.000000	1.000000	1.0	1.000000	1.000000

Figure 4

Looking at how the data correlates with the target column, the target column is 'event_offer completed', I can decide which columns to keep and which columns to discard

	age	income	reward	difficulty	duration	gender_n	web	email	mobile	social	event_offer completed
age	1.000000	0.307404	-0.001495	0.000954	0.003021	-0.141609	0.001358	NaN	-0.002464	0.000440	0.111676
income	0.307404	1.000000	0.000001	-0.000410	-0.000588	-0.216776	-0.003635	NaN	0.001800	0.003278	0.239321
reward	-0.001495	0.000001	1.000000	0.062863	-0.457199	-0.001371	-0.602121	NaN	0.032508	0.323159	-0.130458
difficulty	0.000954	-0.000410	0.062863	1.000000	0.642819	-0.003196	-0.030027	NaN	-0.883902	-0.356764	-0.096435
duration	0.003021	-0.000588	-0.457199	0.642819	1.000000	0.000124	0.055224	NaN	-0.583582	-0.325540	0.038337
gender_n	-0.141609	-0.216776	-0.001371	-0.003196	0.000124	1.000000	-0.000872	NaN	0.003667	0.003422	-0.154269
web	0.001358	-0.003635	-0.602121	-0.030027	0.055224	-0.000872	1.000000	NaN	-0.144437	-0.294620	0.052766
email	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
mobile	-0.002464	0.001800	0.032508	-0.883902	-0.583582	0.003667	-0.144437	NaN	1.000000	0.490247	0.083343
social	0.000440	0.003278	0.323159	-0.356764	-0.325540	0.003422	-0.294620	NaN	0.490247	1.000000	0.053393
event_offer completed	0.111676	0.239321	-0.130458	-0.096435	0.038337	-0.154269	0.052766	NaN	0.083343	0.053393	1.000000

Figure 5 (correlation coefficient matrix)

Exploratory Visualization

The chart below shows a gender distribution of the customers at StarBucks.

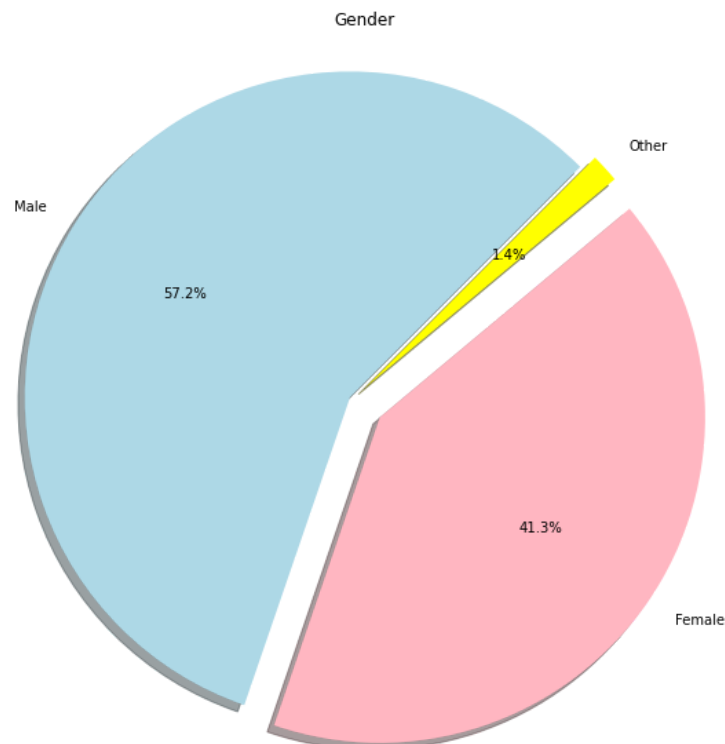


Figure 6 (gender pie chart)

The majority of the customers are male, meaning that most of the transactions will be done by the male customers, if the offers were given to customers at random chances are that the male customers would get the majority of the offers.

The following plot shows which type of offer is likely to be completed.

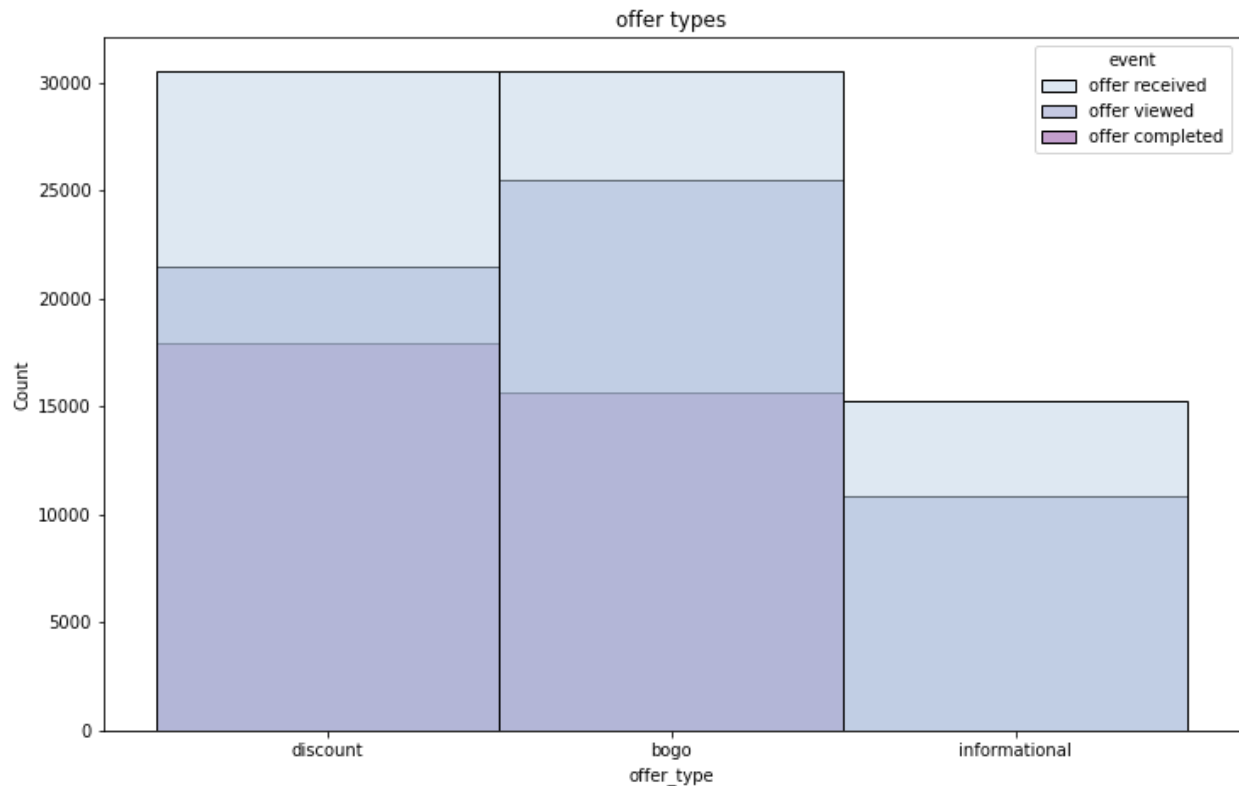


Figure 7

The above plot shows that people are slightly more likely to complete a discount offer than a bogo type offer. Very few people view their discount offer but still complete them still. While the bogo type offers are viewed much more than discount but are not completed as much.

The below scatter plot shows the relationship of offer completion and income of the customer

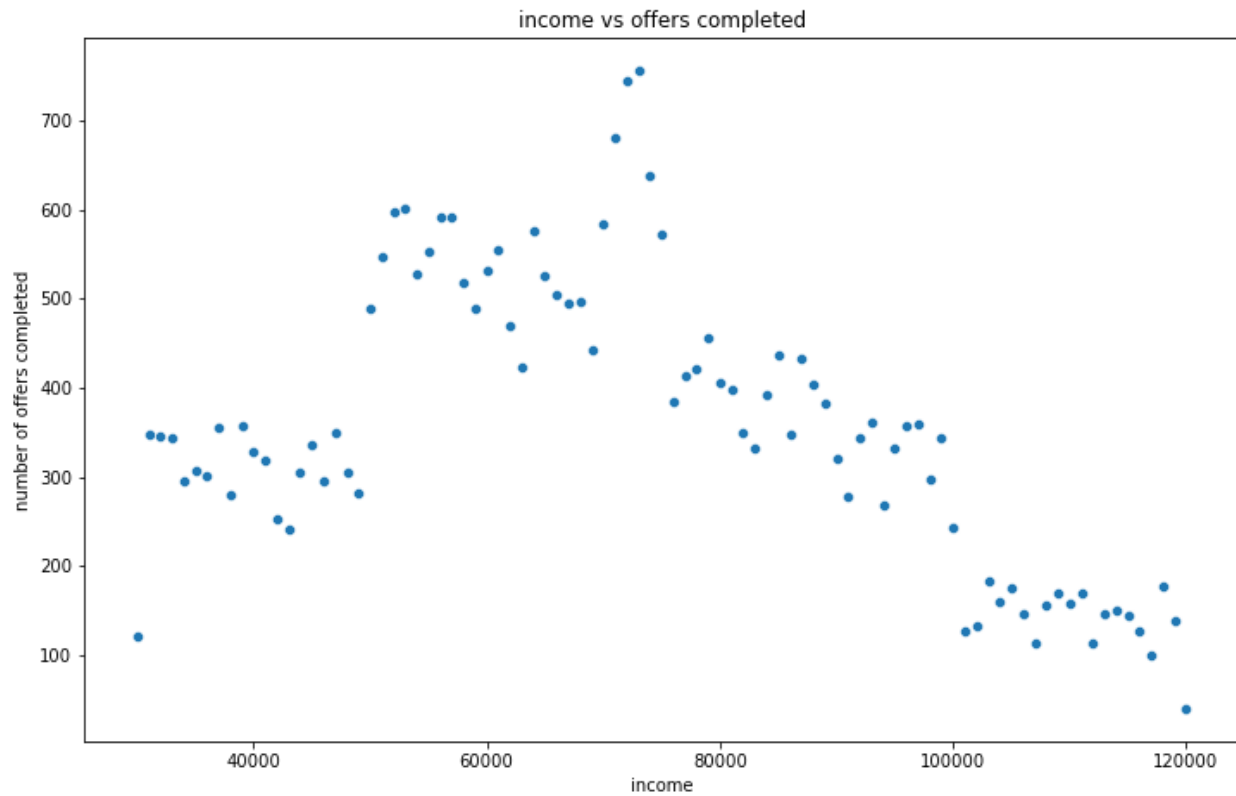


Figure 8

It appears that for customers with high income complete few offers and the higher the income the fewer the offers.

Algorithms and Techniques

The benchmark algorithm is KNN, it is a simple to implement algorithm with simple to adjust hyperparameters. The algorithm computes the distance between data points and input points then classifies the input points as part of the majority points it is nearest to. The KNN uses multiple distance metrics but in this project I will use the euclidean distance metric. The distance metric is as follows

$$d(x, x') = \sqrt{(x_1 - x'_1)^2 + \dots + (x_n - x'_n)^2}$$

The KNN decides the class of a data by using the majority of the nearest neighbors and calculating the probability, using the equation below

$$P(y = J|X = x') = \frac{1}{K} \sum_i (I(y^{(i)} = J))$$

KNN has the following advantages:

- KNN is robust to noisy training data
- It is very effective if the training data is large and it is large in this problem.

KNN has the following disadvantages:

- You need to determine the hyperparameter K (number of neighbors)
- It is not always certain which distance metric produces the best results.
- The computation cost is high because the algorithm computes the distance for the query against all training examples to determine its class.

After processing the StarBucks data, my approach to the problem started to change. I decided to make it into a binary classification problem. The algorithm determines if a customer would complete a given offer or not, to find an offer that a customer would complete you'll have to iterate through the offers and combine the data with customer data and input into the algorithm to get an output.

I used a Decision tree and a Random Forest classifier as well to see how they would do with the given data.

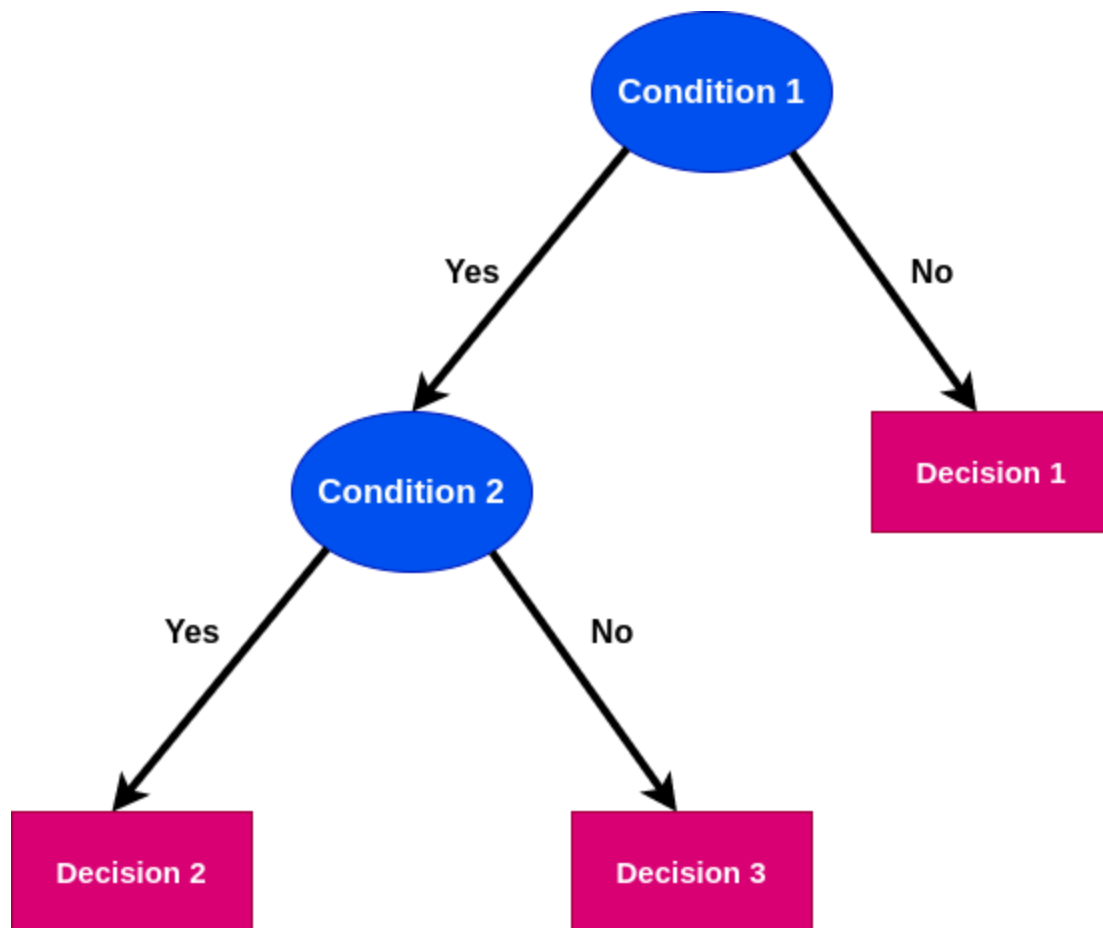


Figure 9 (decision tree depiction)[2]

Figure 9 shows a simplified way on how a decision tree makes decisions. The diagram shows that on each node of the decision tree a condition is formed to separate labels from the dataset until classifications contain features that describe the same thing are made and used to make the final decisions, for instance in this project the final decision should be either 0 or 1. The conditions are set in such a way that entropy is minimized.

$$E = - \sum P_i \log(P_i)$$

Advantages of decision trees:

- It can capture nonlinear relationships: They can be used to classify non-linearly separable data.
- Easy to understand, interpret and visualize.
- It gives us a good idea about the relative importance of attributes.
- Less data preparation needed: In the decision tree, there is no effect by the outsider or missing data in the node of the tree, that's why the decision tree requires fewer data.

- Decision tree is non-parametric: Non-Parametric method is defined as the method in which there are no assumptions about the spatial distribution and the classifier structure.

Disadvantages of decision tree:

- Decision tree for many features: Take more time for training-time complexity to increase as the input increases.
- Decision trees are prone to overfitting, especially when a tree is particularly deep
- It can't be used in big data: If the size of data is too big, then one single tree may grow a lot of nodes which might result in complexity and leads to overfitting.

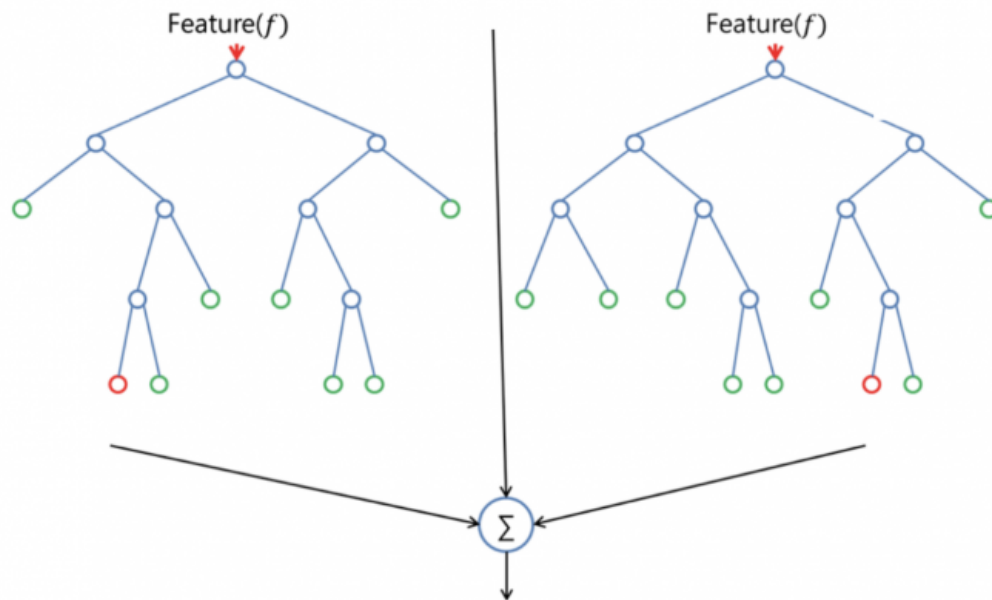


Figure 9 (RandomForest depiction)[3]

A Random forest is similar to a decision tree, as a matter of fact it is made up of decision trees. The random forest algorithm randomly picks rows from the data with replacement, it picks a certain number of rows to make a new dataset. The process by which it makes a dataset is repeated a couple of times and the algorithm ends up having produced multiple datasets. The process of making the datasets is called Bootstrapping. Decision tree algorithms are used on each of the datasets and since it is multiple trees you end up with a forest hence RandomForest. The decision trees do not use all the features but select the features randomly to use. When classifying the algorithm runs the query through all the decision tree models and decides on the class of the query by majority.

Advantages of random forest:

- Random Forest algorithm is less prone to overfitting than Decision Tree and other algorithms
- Random Forest algorithm outputs the importance of features which is a very useful

Disadvantages of random forest:

- Random Forest algorithms may change considerably by a small change in the data.
- Random Forest algorithm computations may grow far more complex compared to other algorithms.

I chose the above algorithm because I needed to classify data and the algorithms are able to do so and do so with simplicity.

Benchmark

The benchmark model was inspired by the Netflix model [\[1\]](#) for recommendation. It uses a KNN algorithm. I have applied it on the given data and I could not get the same average precision, the model has a precision average in the 80% range. The recall I got is an average in the 60% range while netflix seems to have one in the 18% range. The benchmark model has an accuracy of ~69% and an average f1-score of 64%

	precision	recall	f1-score
0	0.63	0.40	0.49
1	0.71	0.86	0.78
accuracy			0.69
macro avg	0.67	0.63	0.64
weighted avg	0.68	0.69	0.67

Methodology

Data Preprocessing

How I processed the data for analysis and use.

I read in the data to my jupyter notebook using pandas.

```
1 portfolio = pd.read_json('portfolio.json', lines=True)
2 profile = pd.read_json('profile.json', lines=True)
3 transcript = pd.read_json('transcript.json', lines=True)
```

Code snippet 1

I merged the data based on common unique columns. The profile and transcript files share the same unique column values which I used to join the two columns

```
1 profile_transcript_df = pd.merge(profile, transcript, left_on='id', right_on='person').drop(columns='person')
```

	gender	age	id	became_member_on	income	event	value	time
0	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	offer received	{'offer id': '2906b810c7d4411798c6938adc9daaa5'}	168
1	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	offer viewed	{'offer id': '2906b810c7d4411798c6938adc9daaa5'}	216
2	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	offer received	{'offer id': '0b1e1539f2cc45b7b9fa7c272da2e1d7'}	336
3	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	offer viewed	{'offer id': '0b1e1539f2cc45b7b9fa7c272da2e1d7'}	348
4	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	transaction	{'amount': 0.35000000000000003}	360

Data snippet 1

I cleaned the column with the name “value” in data snippet 1, so that I can get the unique offer id values in the column to use to join the portfolio file.

After cleaning the column, I only have the values in the dictionary. I removed the none id type of values

I used the following code to clean the data

```
1 def clean_value(value):
2     """
3     This functions gets the value that identifies the offer, thus the unique value of the offer
4     """
5     if 'offer id' in value:
6         return value['offer id']
7     elif 'offer_id' in value:
8         return value['offer_id']
9     else: return np.nan
```

Code snippet 2

```
1 new_value = profile_transcript_df.value.apply(clean_value)
```

```
1 profile_transcript_df['new_value'] = new_value
```

```
1 profile_transcript_df = profile_transcript_df.drop(columns='value')
```

Code snippet 3

Data snippet 2, is the cleaned version of data snippet 1. I have change the column name from “value” to “new_value”

	gender	age	id	became_member_on	income	event	time	new_value
0	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	offer received	168	2906b810c7d4411798c6938adc9daaa5
1	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	offer viewed	216	2906b810c7d4411798c6938adc9daaa5
2	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	offer received	336	0b1e1539f2cc45b7b9fa7c272da2e1d7
3	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN	offer viewed	348	0b1e1539f2cc45b7b9fa7c272da2e1d7

Data snippet 2

During analysis, I saved some of the data into files to save computation speed because they were time intensive.

Since some customers have had more than one transaction at StarBucks, their data appears more than once. I then grouped them by their unique customer id and the offer names which removes the repetitive nature of the data. The code snippet below shows how the data was processed

```
16 #making dummy columns /event_offer completed/event_offer received/event_offer viewed/
17 new_df_with_dummies = pd.get_dummies(df, columns=['event']).dropna()
18
19 #The lines below are trying to get rows that show where an offer has been completed or not
20 mean_df = new_df_with_dummies.groupby(['id_x', 'new_offer_names']).mean().drop(
21     columns=['time', 'event_offer completed',
22             'event_offer received',
23             'event_offer viewed'])
24
25 offers_df = new_df_with_dummies.groupby(['id_x', 'new_offer_names']).sum()[
26     ['event_offer completed',
27      'event_offer received',
28      'event_offer viewed']
29 ]
30
31 formatted_df = mean_df.merge(offers_df, left_on=['id_x', 'new_offer_names'],
32                             right_on=['id_x', 'new_offer_names'])\
33     .reset_index(level=['id_x', 'new_offer_names'])
```

Code snippet 4

Implementation

During the implementation of the algorithms, I loaded the data in the benchmark KNN model after it had been prepared and standardized. The resulting metric outputs were lower than I had hoped. Then I proceeded to implement tree type classification models, the models are Decision tree and Random forest. I made a dictionary with parameters to choose from during hyperparameter tuning to select the best model. The parameters are as follows:

```
36
37 parameters = {
38     'max_depth': [1, 2, 3, 4, 5, 6, 8, 9, 20],
39     'min_samples_leaf': [1, 2, 3, 4, 5, 6, 8, 10, 15, 20, 50, 100],
40     'min_samples_split': [2, 3, 4, 5, 6, 8, 10, 15, 20, 50, 100],
41     'criterion': ['gini', 'entropy']
42 }
43
```

Code snippet 5

The models take in the prepared data as input and use the parameters to fit the best model. The benchmark model is trained on a range of K values ranging from 1-100. The best model was using **K=100 neighbors**. The best decision tree that I trained using grid search had the final parameters as follows:

```
1 best_estimator = chooseSearchGrid(d_tree, 'grid', parameters, X_train, y_train)
DecisionTreeClassifier(criterion='entropy', max_depth=6, min_samples_leaf=50,
                      random_state=10)
```

Code snippet 6

The best decision tree that I train using random search had the following final parameters:

```
1 best_estimator = chooseSearchGrid(d_tree, 'random', parameters, X_train, y_train)
DecisionTreeClassifier(max_depth=5, min_samples_leaf=3, min_samples_split=50,
                      random_state=10)
```

Code snippet 7

The best Random forest model that I trained using grid search had the following final parameters:

```
1 best_estimator = chooseSearchGrid(forest, 'grid', parameters, X_train, y_train)
RandomForestClassifier(max_depth=9, min_samples_leaf=50, random_state=10)
```

Code snippet 8

The best Random forest model that I trained using random search had the following final parameters:

```
1 best_estimator = chooseSearchGrid(forest, 'random', parameters, X_train, y_train)
RandomForestClassifier(criterion='entropy', max_depth=20, min_samples_leaf=8,
                      min_samples_split=100, random_state=10)
```

Code snippet 9

The models are evaluated on precision, recall and f1-score

Refinement

The training Algorithms I used in this project were set in such a way that the best hyperparameters were chosen from the get go of training. There is no initial model for any of the algorithms, only the best model chosen from the provided hyperparameters. The KNN model was trained for different numbers of K values and the best one was chosen. The Random forest and Decision tree classifiers had the same parameters and the best were chosen to be used for the model.

Results

Model Evaluation and Validation

During the training of the algorithms a test data set was used to check for the models' accuracy, precision, recall and f1-score. I compare these evaluation metrics from the different models and then decide on one with the highest of these metrics and f1-score being the most important of the metrics to look out for. All the models I have trained have the same f1 score.

Justification

```
1 best_estimator = chooseSearchGrid(d_tree, 'random', parameters, X_train, y_train)
```

```
DecisionTreeClassifier(criterion='entropy', max_depth=6, min_samples_leaf=3,  
                        min_samples_split=6, random_state=10)
```

```
1 prediction=best_estimator.predict(X_test)
```

```
1 print(classification_report(y_test, prediction))
```

	precision	recall	f1-score	support
0	0.63	0.43	0.51	6519
1	0.72	0.85	0.78	11154
accuracy			0.69	17673
macro avg	0.67	0.64	0.64	17673
weighted avg	0.68	0.69	0.68	17673

Code snippet 10

I chose this DecisionTree model as my final model. This model has a slightly better precision on non completed offers (0) this model will do ever so slightly better than other models. I chose this model based on recall of uncompleted offers represented by 0. The recall is 43% the other models I have trained have a recall just less than 43%. This decision tree model is slightly better than the benchmark model

Conclusion

The insights I have gained from the exploratory data analysis, can help in understanding the distribution of the customer demographic. I was able to find the surprising relationship between income and offers completed. I found that when a customer has high income the customer will likely complete fewer offers. I am able to tell that discount type offers are more likely to be completed than bogo type offers. The models trained have not given impressive results. I would need to go back to the data and do more analysis on it, I would need to try other techniques such as principal component analysis. I would need to do more research on similar projects to try and figure out the feature engineering techniques I could employ for this project

Reference

1. Molina, Leidy Esperanza. "Recommendation System for Netflix."
Recommendation System for Netflix, vol. 1, no. 1, 2018, p. 22.
www.cs.vu.nl,
<https://www.cs.vu.nl/~sbhulai/papers/paper-fernandez.pdf>. Accessed
06 March 2022.
2. Abhijit Roy. "A Dive Into Decision Trees".
www.towardsdatascience.com,
<https://towardsdatascience.com/a-dive-into-decision-trees-a128923c9298#:~:text=A%20decision%20tree%20is%20a,dataset%20to%20the%20fullest%20purity>. Accessed 06 March 2022.
3. Niklas Donges. "A Complete Guide to the Random Forest Algorithm".
www.builton.com,
<https://builton.com/data-science/random-forest-algorithm> . Accessed
06 March 2022