

CIND 820 Capstone Project Final Report: In-Vehicle Coupon Recommendation

Sucheta Sikdar (ID: 501074722)

Supervisor: Dr. Ceni Babaoglu

Toronto Metropolitan University

Date of Submission: 2024-11-30



Table of Contents

Abstract.....	3
Introduction	7
Literature Review	9
Data Description.....	16
Dimensionality Reduction	41
Classification Algorithms and Cross Validation.....	49
Comparison and Analysis of results	56
Methodology	58
Limitations / Challenges / Continuity	60
GitHub Repository.....	61
References.....	61

Abstract

Coupons are a great way for customers to save money on their purchases. It is also a great way for business to attract customers. Coupons create a win-win situation for both companies and customers and hence by offering the correct coupon to users can lead to users to become frequent customers (Niralidedaniya, 2023). If businesses can find the right customers who will use their coupons, then it will help businesses. It will also be interesting to predict what type of coupon will be accepted by a customer based on various attributes about the customer. Luckily, there exists a publicly available dataset called In-Vehicle Coupon Recommendation at UCI Machine Learning Repository from 2020 which describes different driving scenarios of multiple clients such as destination, time, coupon, expiration, gender, age, marital status, whether they have children, education, occupation, income, car, the number of times they go to the bar per month, the number of times they go to a coffee shop per month, the number of times that they buy take away food per month, if customer's average expense per person at restaurants is less than 20 dollars a month, if customer's average expense per person at a restaurant is between 20 dollars to 50 dollars per month, driving distance to the restaurant/bar for using the coupon is greater than 15 minutes, driving distance to the restaurant/bar for using the coupon is greater than 25 minutes, whether the restaurant/bar is in the same direction as destination, whether the restaurant/bar is in the opposite direction as destination, whether the coupon is accepted (*UCI Machine Learning Repository*, 2020).

When printing the info of the dataset in VSCode using Python, the following was observed:

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 12684 entries, 0 to 12683
Data columns (total 26 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   destination                           12684 non-null  object
1   passanger                             12684 non-null  object
2   weather                               12684 non-null  object
3   temperature                           12684 non-null  int64
4   time                                  12684 non-null  object
5   coupon                                12684 non-null  object
6   expiration                             12684 non-null  object
7   gender                                12684 non-null  object
8   age                                   12684 non-null  object
9   maritalStatus                         12684 non-null  object
10  has_children                          12684 non-null  int64
11  education                             12684 non-null  object
12  occupation                             12684 non-null  object
13  income                                12684 non-null  object
14  car                                    108 non-null    object
15  Bar                                    12577 non-null  object
16  CoffeeHouse                           12467 non-null  object
17  CarryAway                             12533 non-null  object
18  RestaurantLessThan20                  12554 non-null  object
19  Restaurant20To50                      12495 non-null  object
20  toCoupon_GEQ5min                      12684 non-null  int64
21  toCoupon_GEQ15min                     12684 non-null  int64
22  toCoupon_GEQ25min                     12684 non-null  int64
23  direction_same                        12684 non-null  int64
24  direction_opp                         12684 non-null  int64
25  Y                                       12684 non-null  int64
dtypes: int64(8), object(18)
memory usage: 2.5+ MB
None
PS C:\Users\suc\OneDrive\Desktop\Sucheta_Sikdar_CIND_820>

```

Figure 1: Info on the In-Vehicle Coupon Recommendation Dataset

As seen in Figure 1, the dataset has 12684 records. The dataset has 25 features and 1 target column. Some of the data like marital status and gender are categorical and some of the data like age and temperature are numerical. Using code provided by Niralidedaniya (2023), it was found that the target classes are partially balanced. If the target classes were highly unbalanced, then

this dataset could not be used because the results of supervised learning algorithms used to make predictions would skew towards the class with the class with higher percentage of records.

```
Accepted coupon: 7210 56.843 %  
Rejected coupon: 5474 43.157 %  
PS C:\Users\suche\OneDrive\Desktop\Sucheta_Sikdar_CIND_820>
```

Figure 2: Distribution of Target Classes

It should also be noted that this dataset comes with many missing values. Hence, this dataset requires preprocessing before it can be analyzed with machine learning algorithms. Using the code provided by Niralidedaniya (2023), we can see the features which have missing values.

	missing_count	missing_percentage
destination	0	0.000000
passanger	0	0.000000
weather	0	0.000000
temperature	0	0.000000
time	0	0.000000
coupon	0	0.000000
expiration	0	0.000000
gender	0	0.000000
age	0	0.000000
maritalStatus	0	0.000000
has_children	0	0.000000
education	0	0.000000
occupation	0	0.000000
income	0	0.000000
car	12576	99.148534
Bar	107	0.843582
CoffeeHouse	217	1.710817
CarryAway	151	1.190476
RestaurantLessThan20	130	1.024913
Restaurant20To50	189	1.490066
toCoupon_GEQ5min	0	0.000000
toCoupon_GEQ15min	0	0.000000
toCoupon_GEQ25min	0	0.000000
direction_same	0	0.000000
direction_opp	0	0.000000
Y	0	0.000000

Figure 3: Distribution of missing values in the dataset

The objective of this project will be to:

- Find the best predictive classification algorithm for the In-Vehicle Coupon Recommendation dataset (2020) after evaluation of various supervised learning classification algorithms introduced to us in *CMTH 642 – Data Analytics: Advanced Methods* like Random Forest, Decision Tree, Logistic Regression, Naïve Bayes, k-Nearest Neighbours (k-NN) on the dataset.
- Using a correlation matrix find out which attributes are highly correlated to the target of the customer accepting or rejecting a coupon.
- Find whether we can attain a dataset with fewer dimensions using these 3 methods: Stepwise Regression, Forward Feature Selection and Backward Feature Elimination methods learnt by us in *CMTH 642 – Data Analytics: Advanced Methods*.
- Find the limitations of this dataset.

Data analysis for this project will be done using Python. Pandas and numpy libraries will be used. Seaborn and Matplotlib libraries will be used for visualizations. Supervised learning techniques will require the scikit-learn library (*Supervised Learning*, n.d.). Evaluation of various classification algorithms will be done by comparing evaluation metrics like the Accuracy and Area under the Curve (AUC) of each algorithm.

Introduction

Coupons are a great way for customers to save money on their purchases and feel special that they are getting a discount. Businesses can attract customers with coupons. If businesses can find the right customers who will use their coupons, then it will help businesses survive and grow. Businesses can then retain existing customers and attract new customers. If businesses know beforehand as to which customers to target for their coupons, then it will save them both money and effort in marketing and sending coupons (Ahmed et al., 2024).

It will also be interesting to predict what type of coupon will be accepted by a customer based on various attributes about the customer. There exists a publicly available dataset called In-Vehicle Coupon Recommendation at UCI Machine Learning Repository from 2020 which describes different driving scenarios of multiple clients and whether the coupon is accepted (*UCI Machine Learning Repository*, 2020).

This project will use the knowledge that has been gained in previous courses of the Data Analytics, Big Data, and Predictive Analytics Certificate Program taught at Toronto Metropolitan University. What has been learned can be implemented in this project.

The objective of this project will be to:

- Find the best predictive classification algorithm for the In-Vehicle Coupon Recommendation dataset (2020) after evaluation of various supervised learning classification algorithms introduced to us in *CMTH 642 – Data Analytics: Advanced Methods* like Random Forest, Decision Tree, Logistic Regression, Naïve Bayes, k-Nearest Neighbours (k-NN) on the dataset. What has been learned will be implemented.

- Using a correlation matrix find out which attributes are highly correlated to the target of the customer accepting or rejecting a coupon. A visual of a correlation matrix will be very effective. **Since most of the features are categorical, the discrete values will be converted to corresponding numerical values using encoding.**
- Find whether we can attain a dataset with fewer dimensions using these 3 methods: Stepwise Regression, Forward Feature Selection and Backward Feature Elimination methods learnt by us in *CMTH 642 – Data Analytics: Advanced Methods*. This will help eliminate some noise. **Dimensionality reduction is the major contribution of this project to this dataset.**
- Find the limitations of this dataset. The flaws of this dataset need to be explored.

There is some existing and related work done on this topic already. Exploratory data analysis to get a better picture of our data is also necessary.

Literature Review

The original data was collected by Wang et al. (2017) using Amazon Mechanical Turk about users interacting with a mobile recommendation system. They used rule set models that are used for classification and decision making to understand and predict consumers' response to coupons in different contexts (Wang et al., 2017). Wang et al. (2017) used Bayesian approach and Disjunctive Normal Form classifiers. An example is that if a combination of a subset of input is met, then the output is satisfied. Wang et al. (2017) compared performance in accuracy by calculating and comparing Area under the ROC curve for their deduced final two Bayesian Rule Sets with other machine learning methods including different types of decision tree algorithms. Wang et al. (2017) compared with C4.5, CART, Lasso, RIPPER and TopK. Unlike other machine learning algorithms, a Bayesian approach looks at previous choices. Wang et al. (2017) deduced that that their Bayesian approach had competitive performance. Their methods are more complex compared to the simpler supervised learning classification algorithms introduced to us in *CMTH 642 – Data Analytics: Advanced Methods*. It is worth checking how this dataset performs against simpler supervised learning classification algorithms introduced to us in *CMTH 642 – Data Analytics: Advanced Methods*. It is possible that using simpler supervised learning classification algorithms like Random Forest, Decision Tree, Logistic Regression, Naïve Bayes, k-Nearest Neighbours (k-NN) might yield similar competitive results too.

Same dataset was used by Niralidedaniya (2023). A lot of the data understanding, Data preparation and Exploratory Data Analysis on the dataset was already done by Niralidedaniya (2023). Niralidedaniya (2023) experimented with below thirteen machine learning classification models with hyper parameter tuning:

1. Logistic Regression

2. K-Nearest Neighbor
3. Decision Tree
4. Support Vector Classification with rbf kernel
5. LinearSVC
6. Gaussian Naive Bayes
7. Random Forest
8. GBDT
9. Bagging Classifier
10. AdaBoost Classifier
11. Gradient Boosting Classifier
12. ExtraTrees Classifier
13. HistGradientBoosting Classifier

Niralidedaniya (2023) compared the test log loss and test AUC score of each of the thirteen ML models with different encoded data matrices to find the best model and best encoding method with their best hyper parameter. Niralidedaniya (2023) found that Ordinal Encoding and One Hot Encoding perform best than other encoding techniques. It was observed that XGB Classifier, Support Vector Classification, Hist Gradient Boosting, and Random Forest Classifier models perform best than other models. Bagging, AdaBoost, Gradient Boosting, and ExtraTrees Classifier models perform best but they are overfitted (Niralidedaniya, 2023). Niralidedaniya (2023) stated that their Training AUC Score is very high with a value of 1, and a training log loss is almost 0. Basic classification models like Logistic Regression, K-Nearest Neighbor, Decision Tree, LinearSVC, and Gaussian Naive Bayes models didn't do well for this problem (Niralidedaniya, 2023). Niralidedaniya (2023) took classification up a notch by the use of

stacking classifier on the four models: XGB Classifier, Support Vector Classification, Hist Gradient Boosting, and Random Forest Classifier with their best parameters. It would still be interesting to see and compare how the basic classification algorithms introduced to us in *CMTH 642 – Data Analytics: Advanced Methods* like Random Forest, Decision Tree, Logistic Regression, Naïve Bayes, k-Nearest Neighbours (k-NN) perform on the dataset. A generated visual correlation matrix will be insightful. A dataset with fewer dimensions can be attained using these 3 methods: Stepwise Regression, Forward Feature Selection and Backward Feature Elimination methods learnt by us in *CMTH 642 – Data Analytics: Advanced Methods*. It will also be useful to explore the limitations of this dataset.

Depari et al. (2022) used the same dataset. Depari et al. (2022) used RapidMiner 9.10.001 tool in predicting customer's responses to in-vehicle coupon recommendations. Depari et al. (2022) compared the accuracy percentage, class precision, and execution time for three algorithms (Random Forest, Naïve Bayes and Decision Tree) on the dataset after doing a detailed descriptive analysis of the data. The descriptive analysis was very insightful about the dataset. For example, Depari et al. (2022) found that the data contained mostly married females who like to travel alone on a sunny day around 6 PM. Most of them have attended college, yet didn't graduate (Depari et al., 2022). For those who have an occupation, it states that most of them earn an income of around \$25000 - \$37499 (Depari et al., 2022). It was also mentioned that the destination is mostly the No Urgent Place such as Coffee House, which provides a coupon that expires in one day (Depari et al., 2022). It will be interesting to find more inferences from the data with the help of plots for each feature. This will help with finding the limitations of the dataset provided. Depari et al. (2022) found that predictive analytics results showed that random forest achieved the highest accuracy with 77.65% overall accuracy percentage, yet required the

most time to process. However, the decision tree algorithm acquired the highest confidence level of 0.750 for prescriptive analysis (Depari et al., 2022). Prescriptive analytics is data analysis that is used to explain and determine what is the next action plan and why should they do it (Depari et al., 2022). It will be interesting to compare with 2 more algorithms: Logistic Regression and k-Nearest Neighbours (k-NN) used on the dataset. There was no dimension reduction done by Depari et al. (2022).

This dataset was also used by Atiq et al. (2022). This dataset comes with missing values. This could cause problems in the prediction analysis. To circumvent this problem Atiq et al. (2022) analysed the impact of four different imputation techniques (Frequent value, mean, KNN, MICE) to replace the missing values and use them to create prediction models. Atiq et al. (2022) then applied six classifier algorithms (Naïve Bayes, Deep Learning, Logistic Regression, Decision Tree, Random Forest, and Gradient Boosted Tree). Atiq et al. (2022) found that KNN imputation with Deep Learning classifier gave the most accurate outcome. Atiq et al. (2022) applied SMOTE (Synthetic Minority Oversampling Technique) for oversampling the dataset in order to have positive and negative instances divided into 50% / 50%. This led to a perfectly balanced dataset for target value. Below is a table by Atiq et al. (2022) that is used to show Accuracy of Classifiers on Actual and Oversampled Dataset using various imputation Techniques. **These pre-existing findings will be useful in dealing with the missing values in the dataset.** Also, there was no dimension reduction done by Atiq et al. (2022).

Imputation Technique	Classifier	Accuracy (Without Over Sampling)	Accuracy (With Over Sampling)
Frequent Value	Random Forest	74.2	76.1
	Decision Tree	70.6	69.5
	Logistic Regression	68.433	69.6
	Gradient Boosted Tree	74.338	75.4
	Naive Bayes	66.253	67.1
	Deep Learning	75.73	60.4
Mean Imputation	Random Forest	74.7	73.1
	Decision Tree	68.82	67.3
	Logistic Regression	69.2	67.5
	Gradient Boosted Tree	73.4	73.1
	Naive Bayes	65.9	68.3
	Deep Learning	57.40	100
KNN	Random Forest	74.32	72.3
	Decision Tree	71.91	68.1
	Logistic Regression	68.8	66.2
	Gradient Boosted Tree	73.2	72.4
	Naive Bayes	65.6	64.9
	Deep Learning	57.16	83.27
MICE	Random Forest	74.3	71.9
	Decision Tree	70.2	67.9
	Logistic Regression	66.4	68.1
	Gradient Boosted Tree	73.2	72.9
	Naive Bayes	66.6	65.3
	Deep Learning	100	80.61

Table 1: Accuracy of Classifiers on Actual and Over Sampled Dataset provided by Atiq et al. (2022).

Patil et al. (2019) used a different dataset for E-coupons to predict coupon usage behaviour. It is interesting that to train the gradient boosting classifier, they used 45 “train periods” that simulated the test timing (Patil et al., 2019). When comparing with other algorithms like logistic regression, SVM, random forest, neural networks, Patil et al. (2019) found that gradient boosting was the single best classifier. By having train periods, Patil et al. (2019) made interesting observations such as the more the coupon is viewed, the probability to buy using the coupon code increases. Patil et al. (2019) also observed that customers tend to purchase the same coupon

over and over again. The dataset for in-vehicle coupon response deficient in data over a periodic basis to help uncover such patterns. This is a limitation of the in-vehicle coupon dataset.

Ahmed et al. (2024) used a different dataset, Dunnhumby data for a particular grocery retail company. Ahmed et al. (2024) proposed two different models: one for predicting customer churn and the other for coupon redemption model and both those models used XGBoost Classifier Model. As seen in earlier papers, it is good to use a variety of models for comparison. However, Ahmed et al. used only XGBoost Classifier Model. Ahmed et al. (2024) also stated that their dataset was not large enough and therefore it led to uncertainty about the model performance on a larger dataset in the future.

Given all the existing data analysis that has been done on coupon redemption, it is obvious that it is crucial for businesses to help them grow and survive. When we attain a dataset with fewer dimensions using these 3 methods: Stepwise Regression, Forward Feature Selection and Backward Feature Elimination methods learnt by us in *CMTH 642 – Data Analytics: Advanced Methods*, we can have more concise information to be used for decision making and not have noise. **Dimensionality reduction is the major contribution of this project to this dataset.** A visual correlation matrix will show which attributes are highly correlated to the target of the customer accepting or rejecting a coupon. A visual correlation matrix is a very effective way to display a lot of information. This project also gives an opportunity to compare and find the best predictive classification algorithm for In-Vehicle Coupon Recommendation dataset (2020) after evaluation of various supervised learning classification algorithms introduced to us in *CMTH 642 – Data Analytics: Advanced Methods* like Random Forest, Decision Tree, Logistic Regression, Naïve Bayes, k-Nearest Neighbours (k-NN) on the dataset. There could be so many

other factors that effect customer coupon redemption. Therefore, we need to also explore the limitations of this dataset. Hence, it is worth proceeding with this project.

Data Description

The dataset is in csv format. It is a publicly available dataset called In-Vehicle Coupon

Recommendation at UCI Machine Learning Repository from 2020 which describes different driving scenarios of multiple clients and whether the coupon is accepted (*UCI Machine Learning Repository*, 2020). The dataset is large with 12684 records and 25 features. The below figure provides some information on the dataset:

```
▶ in_vehicle_coupon_data.info()
```

```
↗ <class 'pandas.core.frame.DataFrame'>  
RangeIndex: 12684 entries, 0 to 12683  
Data columns (total 26 columns):  
#   Column                                Non-Null Count  Dtype  
---  ---                                -  
0   destination                          12684 non-null  object  
1   passanger                           12684 non-null  object  
2   weather                             12684 non-null  object  
3   temperature                         12684 non-null  int64  
4   time                               12684 non-null  object  
5   coupon                             12684 non-null  object  
6   expiration                          12684 non-null  object  
7   gender                              12684 non-null  object  
8   age                                 12684 non-null  object  
9   maritalStatus                      12684 non-null  object  
10  has_children                        12684 non-null  int64  
11  education                          12684 non-null  object  
12  occupation                         12684 non-null  object  
13  income                            12684 non-null  object  
14  car                                108 non-null    object  
15  Bar                                12577 non-null  object  
16  CoffeeHouse                       12467 non-null  object  
17  CarryAway                         12533 non-null  object  
18  RestaurantLessThan20              12554 non-null  object  
19  Restaurant20To50                  12495 non-null  object  
20  toCoupon_GEQ5min                  12684 non-null  int64  
21  toCoupon_GEQ15min                 12684 non-null  int64  
22  toCoupon_GEQ25min                 12684 non-null  int64  
23  direction_same                    12684 non-null  int64  
24  direction_opp                     12684 non-null  int64  
25  Y                                  12684 non-null  int64  
dtypes: int64(8), object(18)  
memory usage: 2.5+ MB
```

Figure 4: In-Vehicle Coupon Recommendation dataset information

Feature	Feature Description	Possible Values
destination	Destination	No Urgent Place, Home, Work
passenger	Passenger	Alone, Friend(s), Kid(s), Partner
weather	Weather Type	Sunny, Rainy, Snowy
temperature	Temperature	Numerical value
time	Time	2PM, 10AM, 6PM, 7AM, 10PM
coupon	Coupon type	Restaurant (<20), Coffee House, Carry out & Take away, Bar, Restaurant (20-50)
expiration	Expiration	1d, 2h
gender	Gender	Female, Male
age	Age	below21, 50plus, 36, 41, 31, 26, 46, 21
maritalStatus	Marital Status	Unmarried partner, Single, Married partner, Divorced, Widowed
has_children	Whether they have children	Numerical value – 0 or 1
education	Education	Some college – no degree, Bachelors degree, Associates degree, High School

		Graduate, Graduate degree (Masters or Doctorate), some High School
occupation	Occupation	Unemployed, Architecture & Engineering, Student, Education & Training & Library, Healthcare Support, Healthcare Practitioners & Technical, Sales & Related, Management, Arts Design Entertainment Sports & Media, Computer & Mathematical, Life Physical Social Science, Personal Care & Service, Community & Social Services, Office & Administrative Support, Construction & Extraction, Legal, Retired, Installation Maintenance & Repair, Transportation & Material Moving, Business & Financial, Protective Service,

		Food Preparation & Serving Related, Production Occupations, Building & Grounds Cleaning & Maintenance, Farming Fishing & Forestry
income	Income	37500-49999, 62500-74999, 12500 – 24999, 75000-87499, 50000 – 62499, 25000 – 37499, \$100000 or More, 87500 – 99999, Less than \$12500
car	Car	Scooter and motorcycle, crossover, Mazda5, do not drive, car that is too old to install Onstar
bar	The number of times they go to the bar per month	Never, less 1, 1~3, gt8, 4~8
CoffeeHouse	The number of times they go to a coffee shop per month	Never, less1, 4~8, 1~3, gt8
CarryAway	The number of times that they buy take away food per month	Never, less1, 4~8, 1~3, gt8

RestaurantLessThan20	if customer's average expense per person at restaurants is less than 20 dollars a month	Never, less1, 4~8, 1~3, gt8
Restaurant20To50	if customer's average expense per person at a restaurant is between 20 dollars to 50 dollars per month	Never, less1, 4~8, 1~3, gt8
toCoupon_GEQ5min	driving distance to the restaurant/bar for using the coupon is greater than 5 minutes	Numerical value
toCoupon_GEQ15min	driving distance to the restaurant/bar for using the coupon is greater than 15 minutes	Numerical value
toCoupon_GEQ25min	driving distance to the restaurant/bar for using the coupon is greater than 25 minutes	Numerical value
direction_same	whether the restaurant/bar is in the same direction as destination	Numerical value
direction_opp	whether the restaurant/bar is in the opposite direction as destination,	Numerical value

Y	whether the coupon is accepted. This is the target.	0 or 1. 1 if coupon is accepted. 0 if coupon is rejected.
---	------------------------------------------------------------	-----------------------------------------------------------

Table 2: Feature Description in In Vehicle Coupon Redemption dataset

The below plots for distribution of categorical value features were generated with the help of code by Inyama (2023).

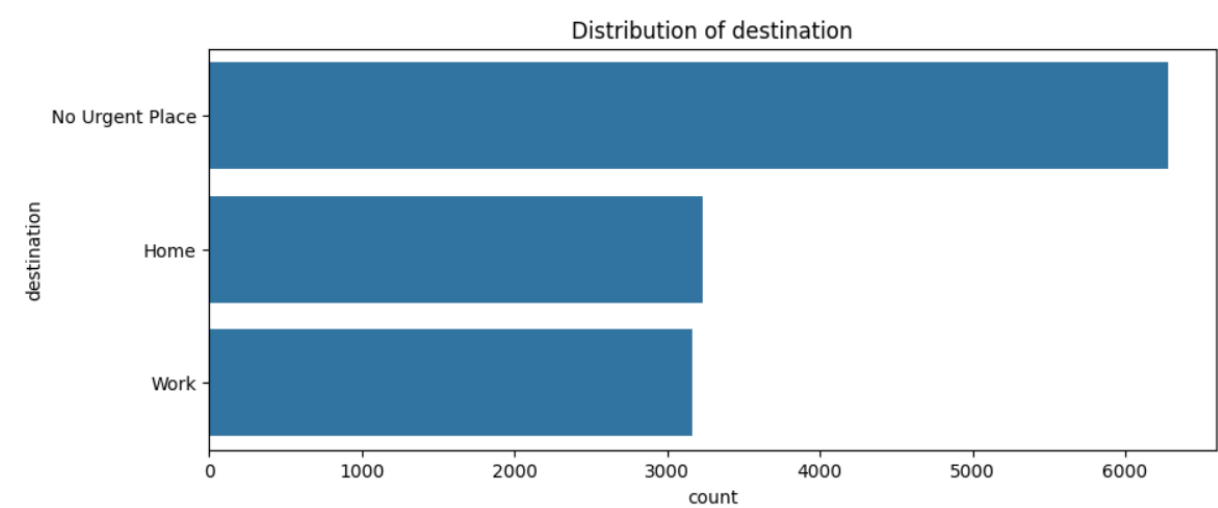


Figure 5: Distribution of destination feature

Observation: Most people’s destination is not an urgent place

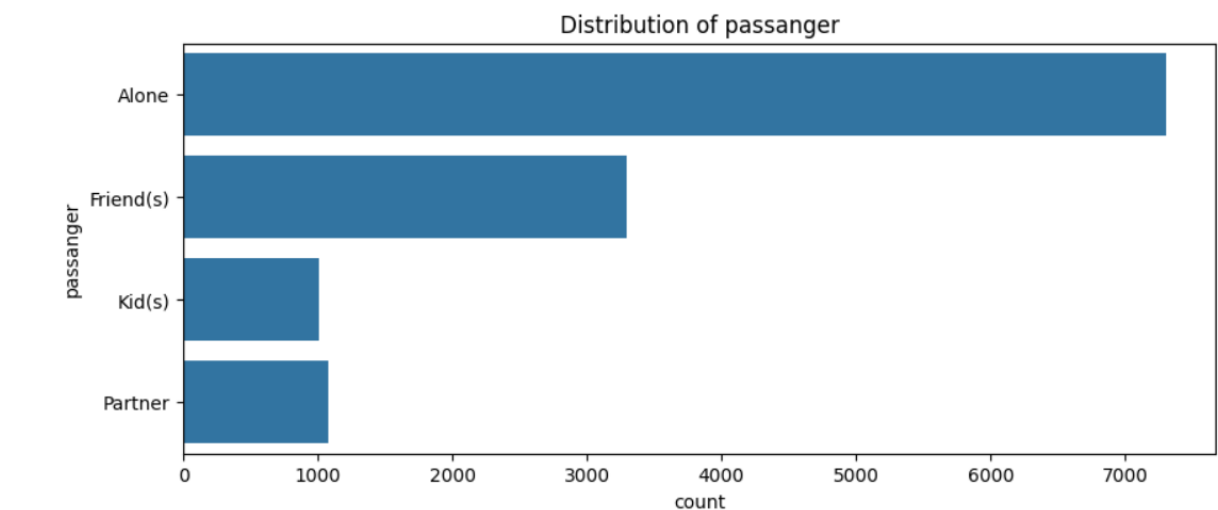


Figure 6: Distribution of passenger feature

Observation: Most people travel alone

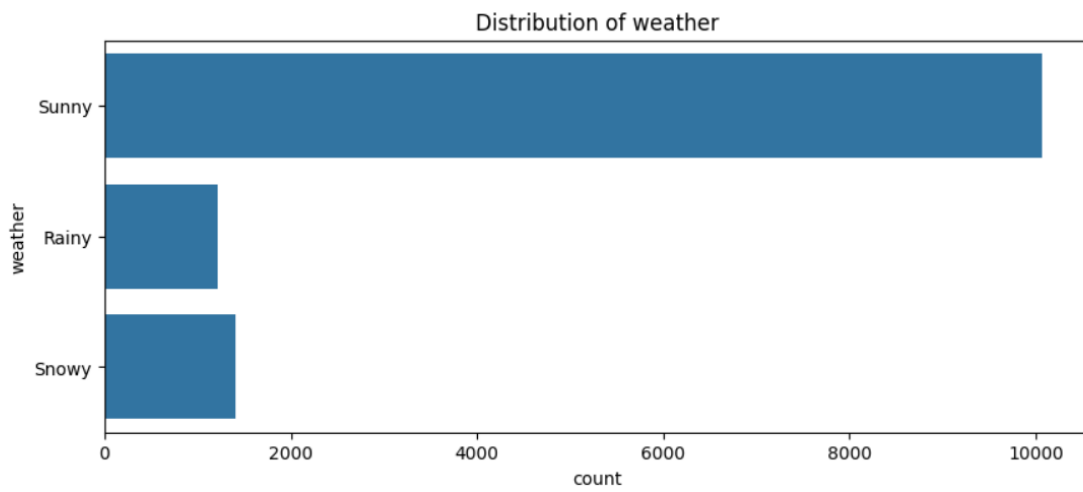


Figure 7: Distribution of weather feature

Observation: Most people travel on a sunny day

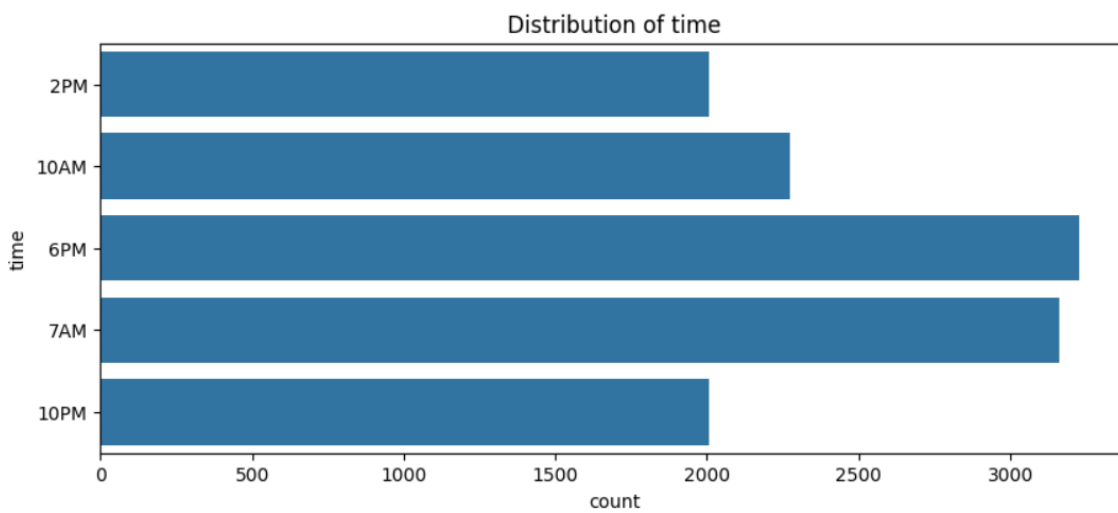


Figure 8: Distribution of time feature

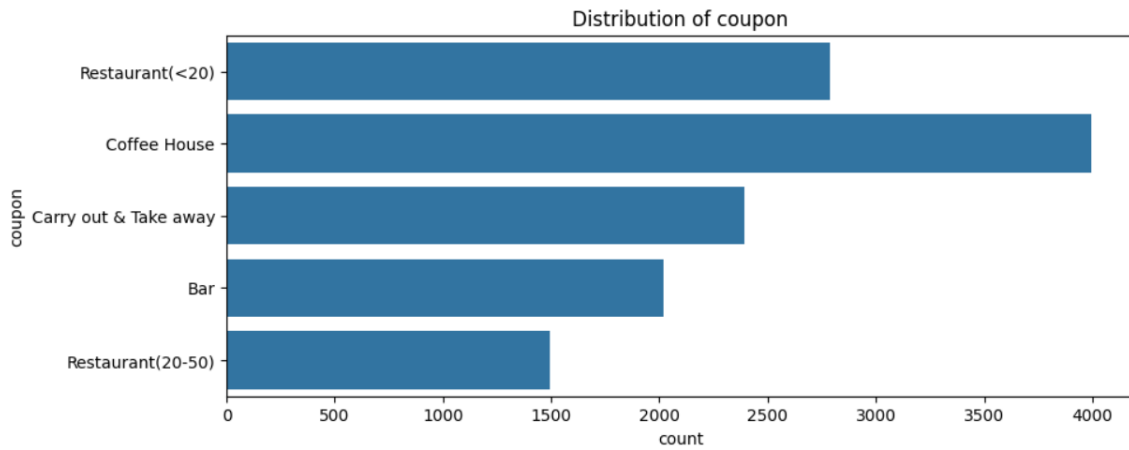


Figure 9: Distribution of coupon feature

Observation: Most people go to Coffee House

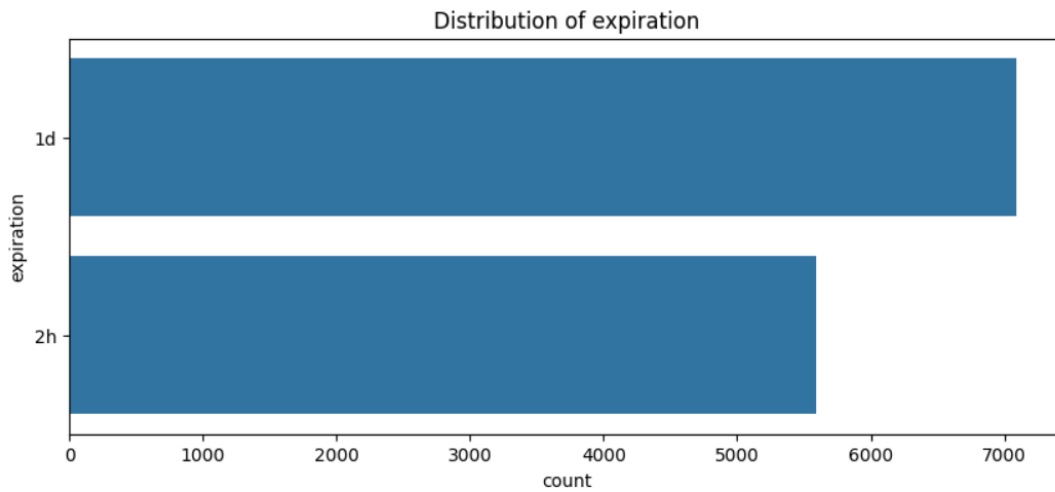


Figure 10: Distribution of expiration feature

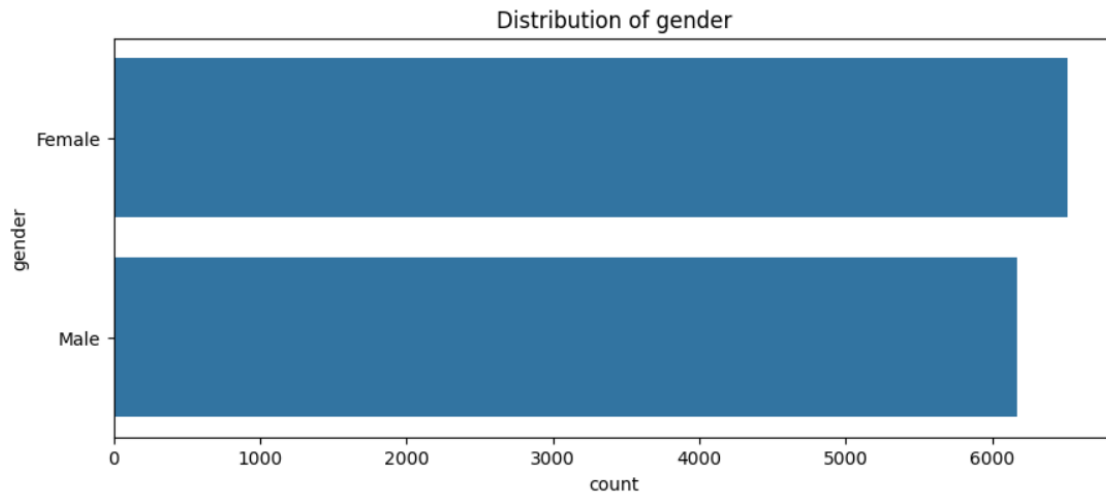


Figure 11: Distribution of gender feature

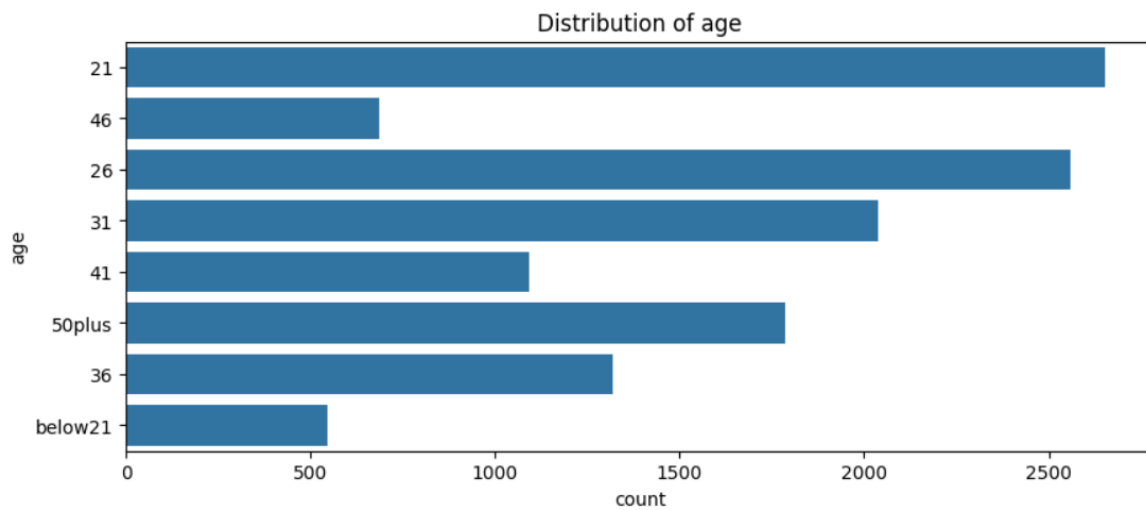


Figure 12: Distribution of age feature

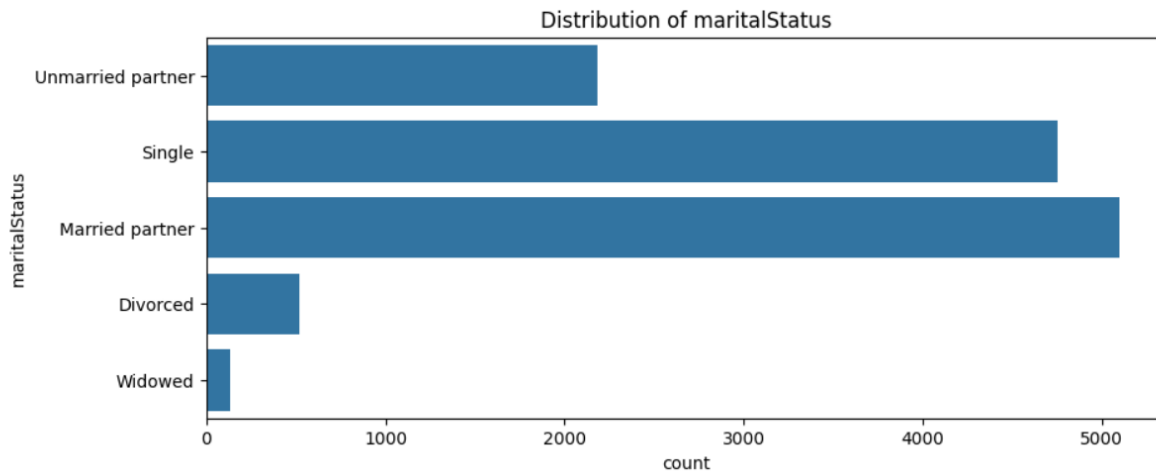


Figure 13: Distribution of marital status feature

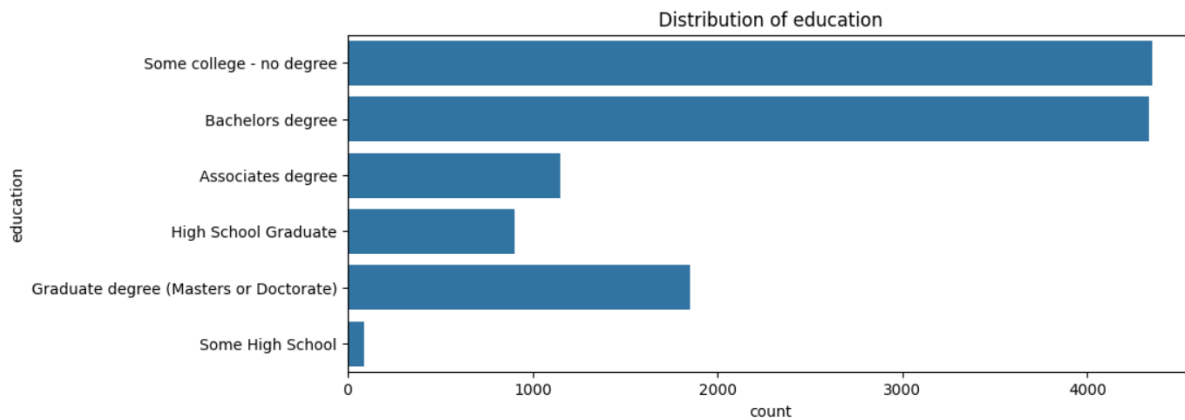


Figure 14: Distribution of education feature

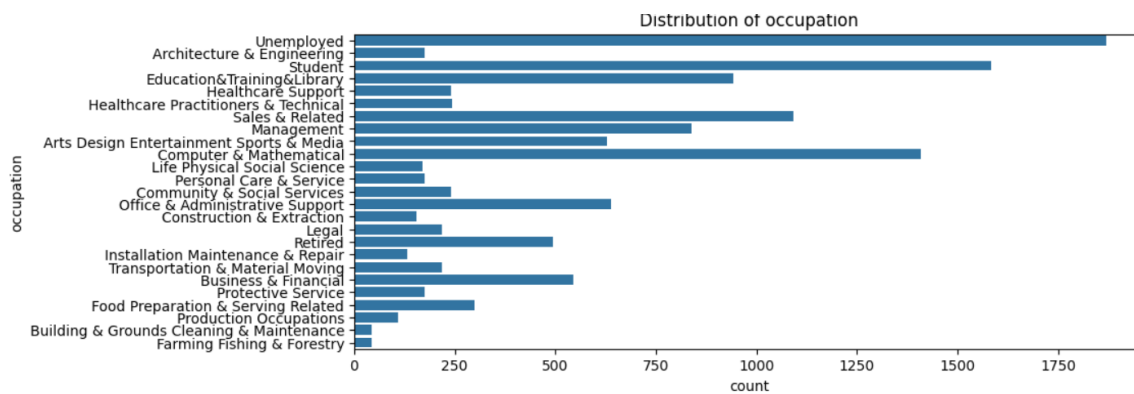


Figure 15: Distribution of occupation feature

Observation: Occupation has too many categories.

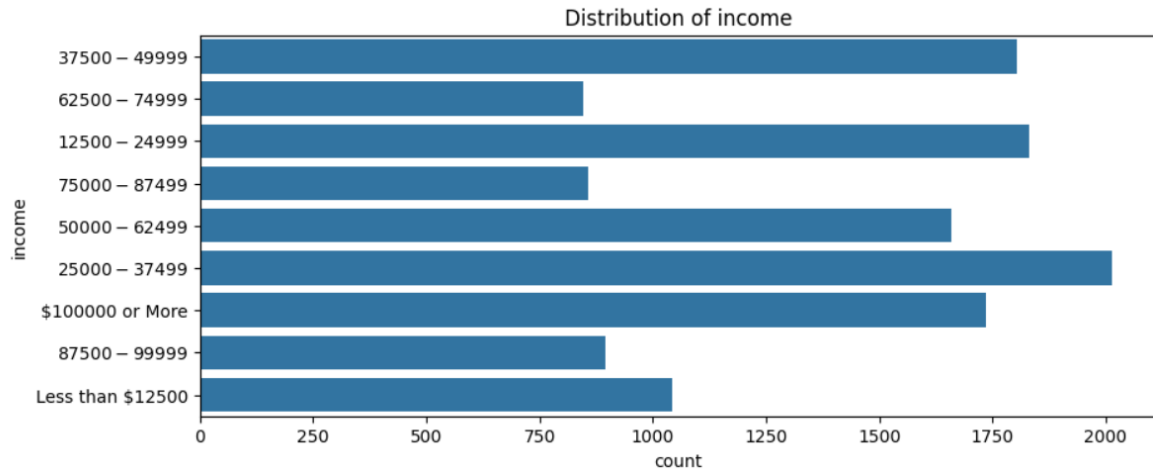


Figure 16: Distribution of income feature

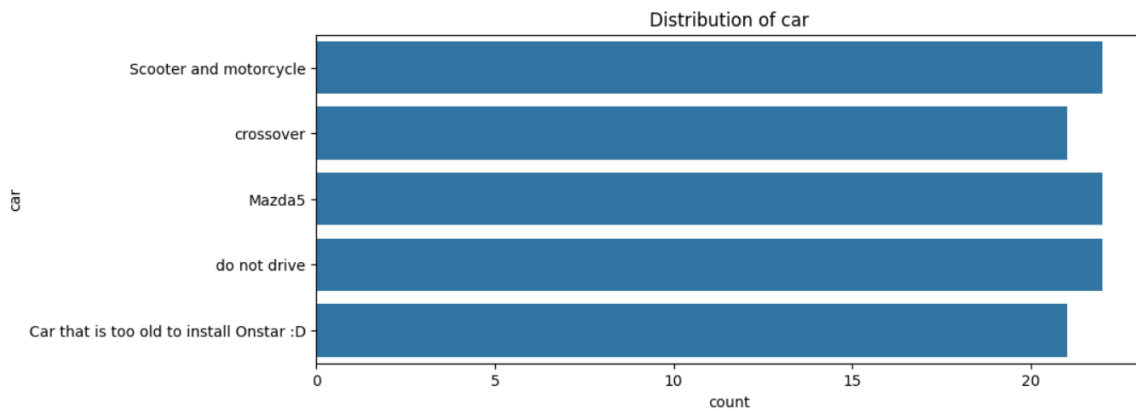


Figure 17: Distribution of car feature

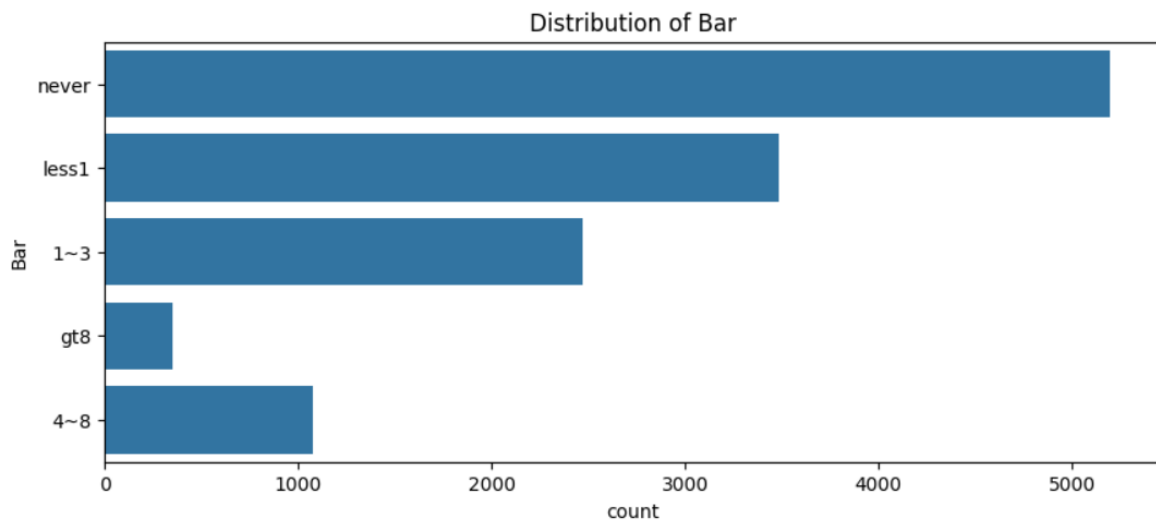


Figure 18: Distribution of bar feature

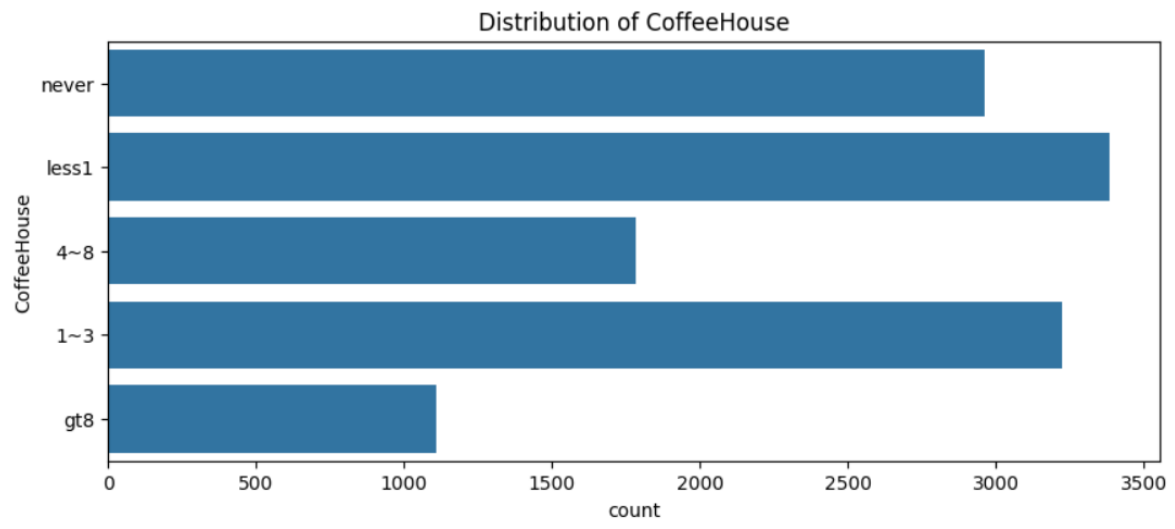


Figure 19: Distribution of Coffee House feature

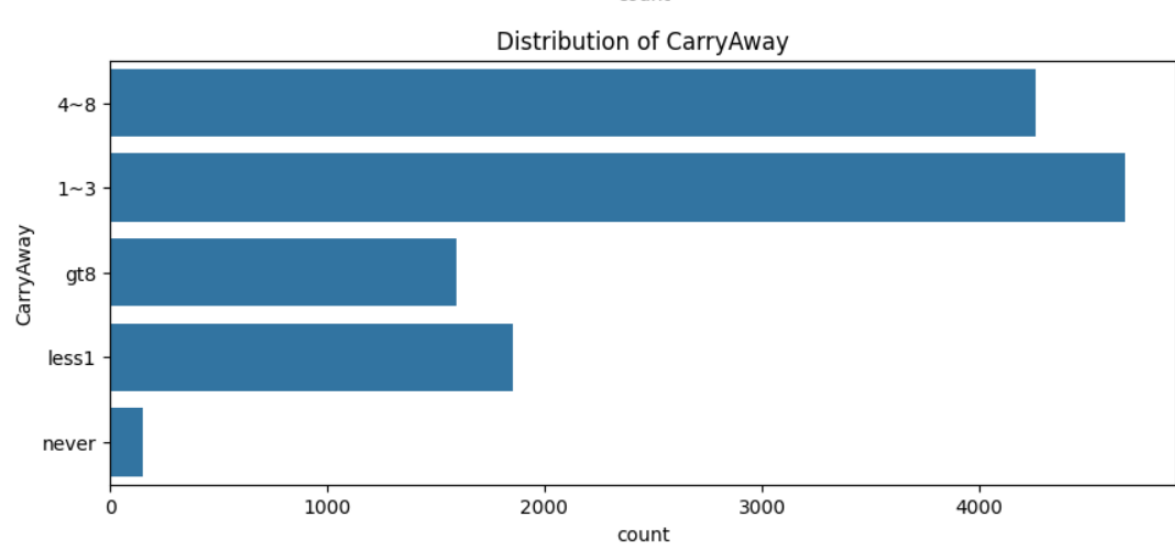


Figure 20: Distribution of Carry Away feature

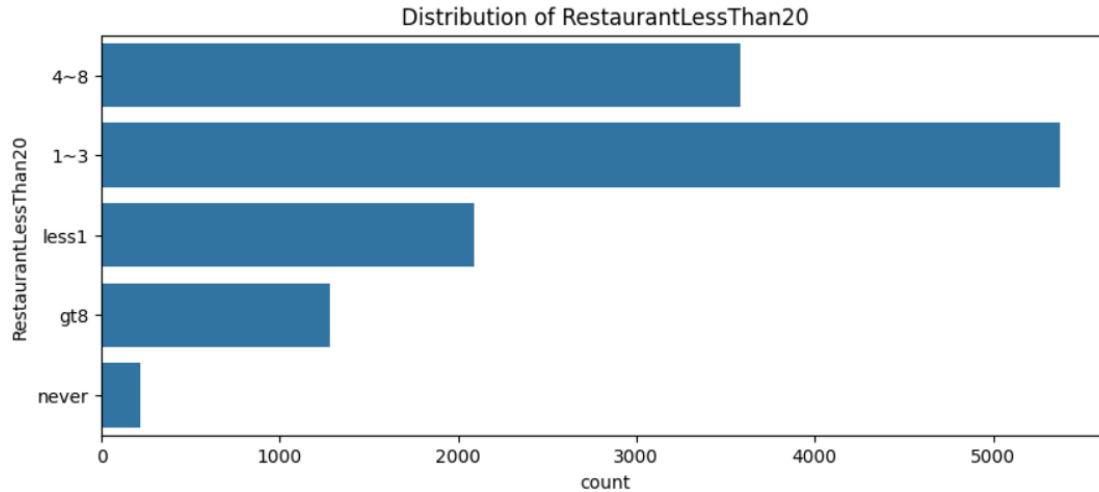


Figure 21: Distribution of Restaurant Less Than 20 feature

Observation: Most popular frequency for RestaurantLessThan20 feature is 1 to 3

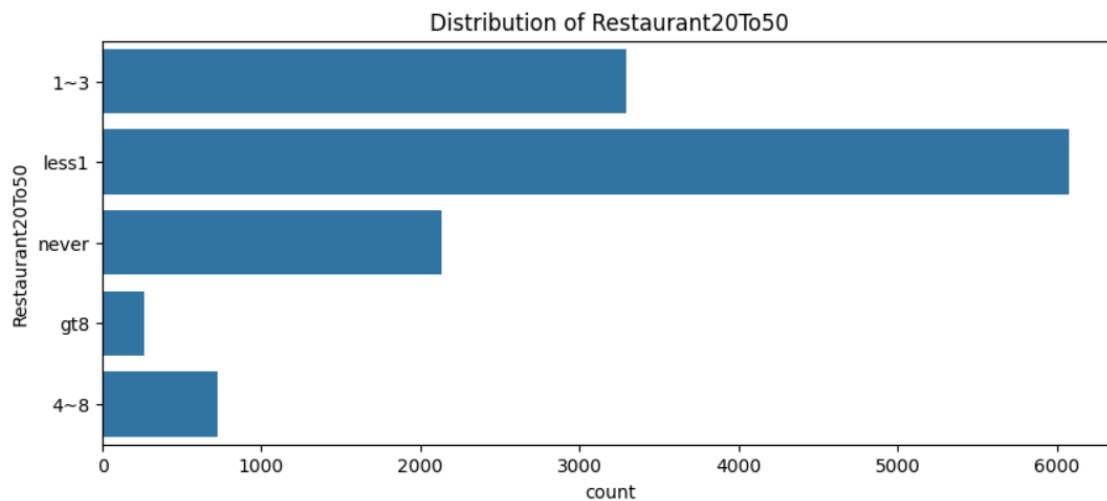


Figure 22: Distribution of Restaurant 20 to 50 feature

Observation: Most popular frequency for Restaurant20To50 feature is less 1

Using code from Inyama (2023), we can see the boxplots of the numerical data in numerical columns.

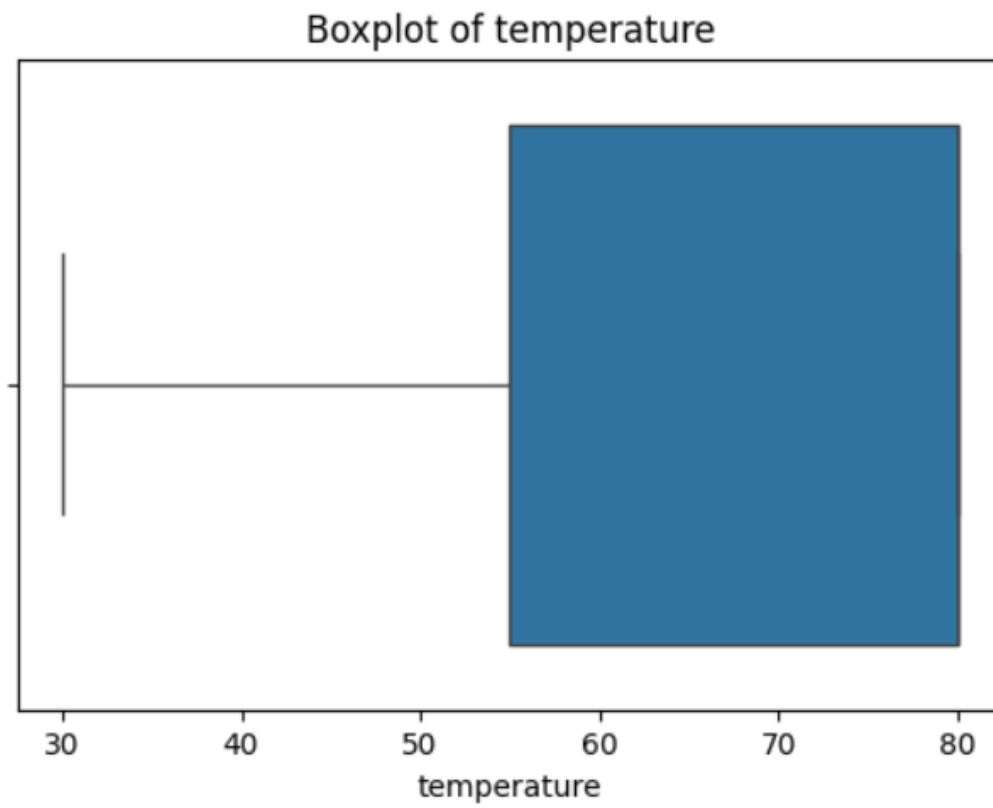


Figure 23: Boxplot of temperature

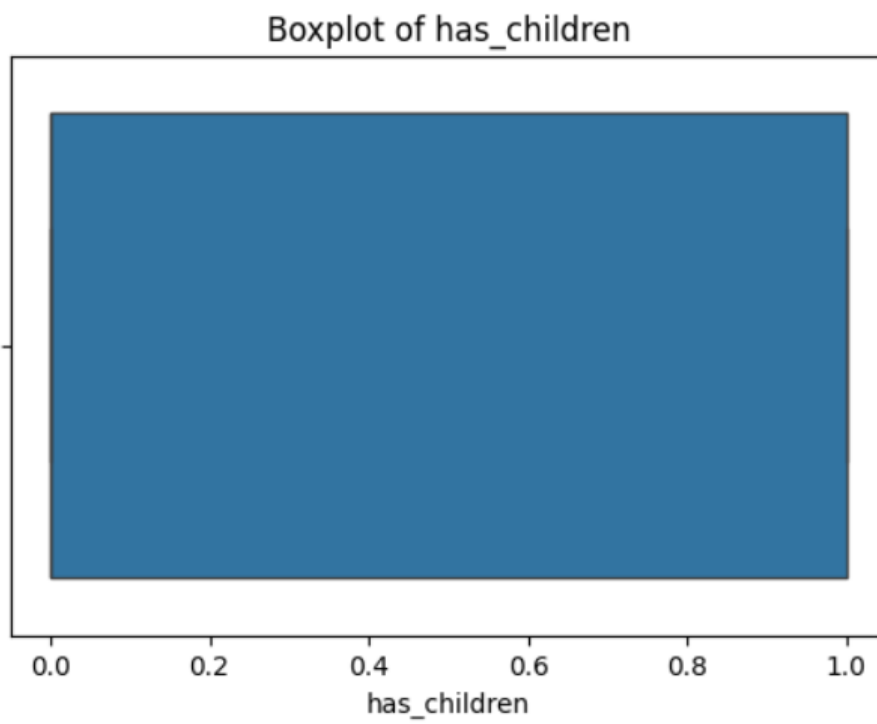


Figure 24: Boxplot of has children

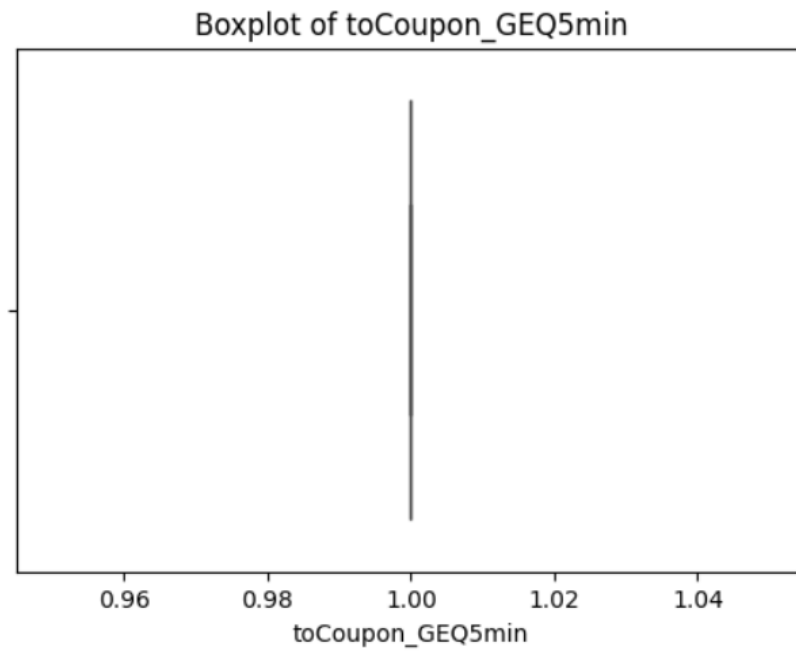


Figure 25: Boxplot of to coupon GEQ 5 min

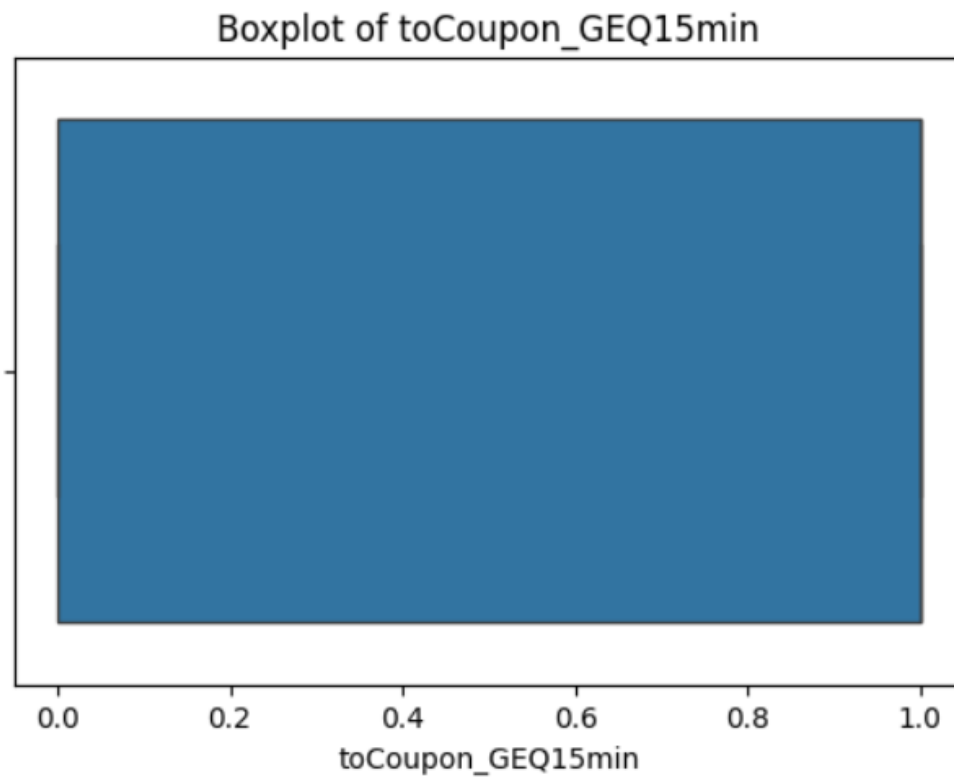


Figure 26: Boxplot of to coupon GEQ 15 min

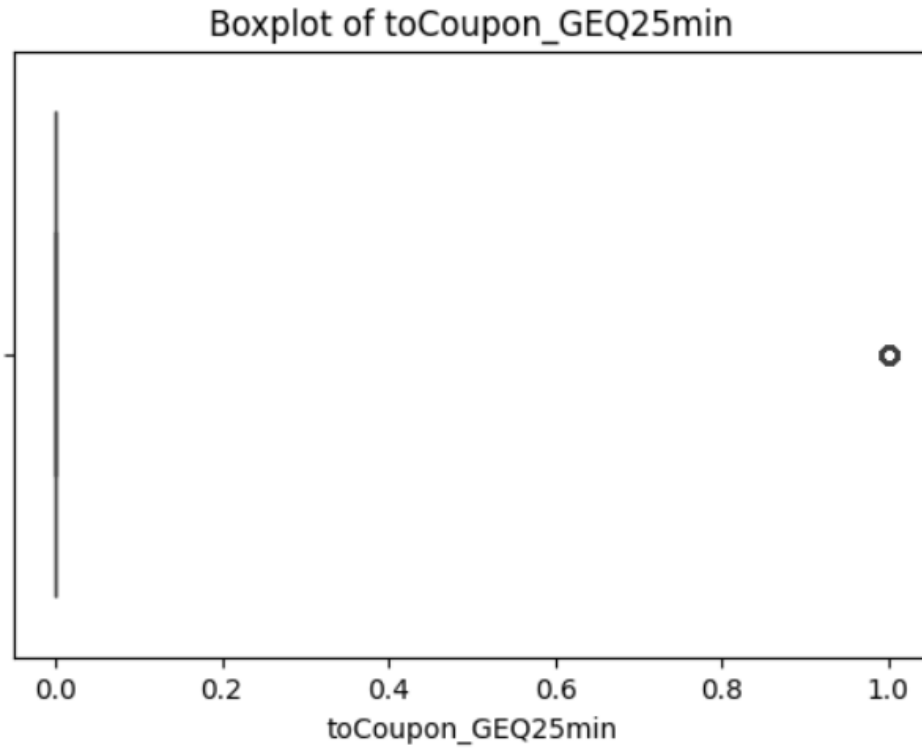


Figure 27: Boxplot of to coupon GEQ 25 min

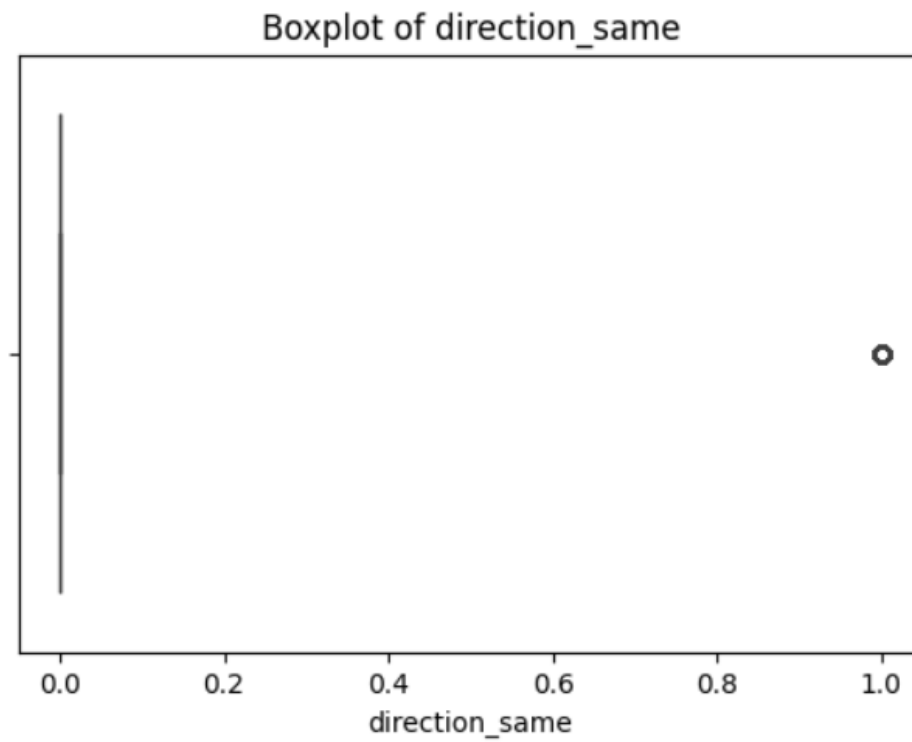


Figure 28: Boxplot of direction same

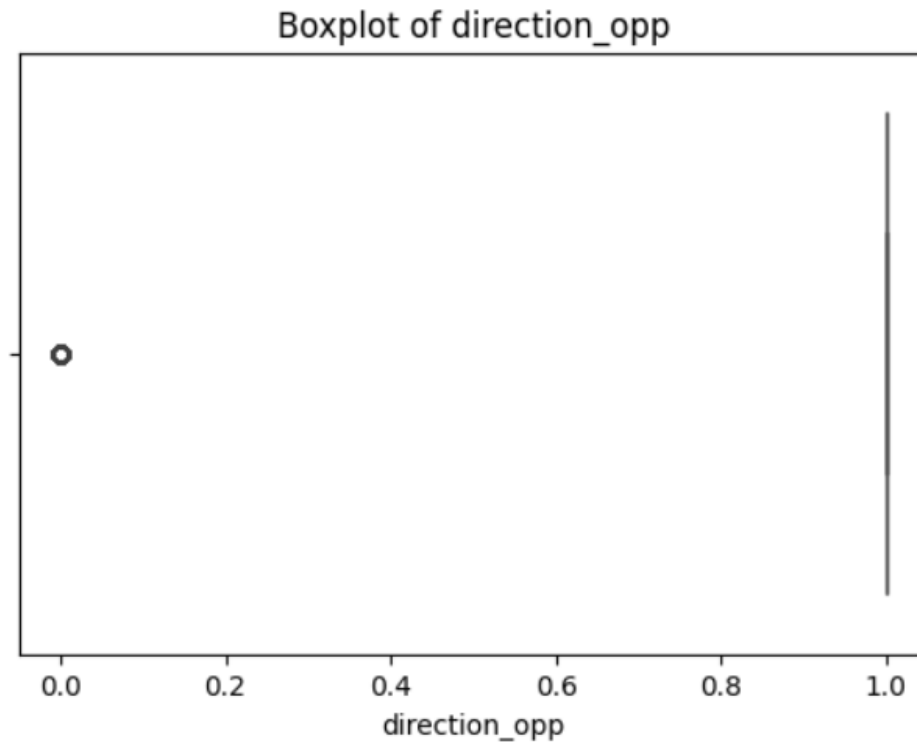


Figure 29: Boxplot of direction opposite

Using code provided by Niralidedaniya (2023), it was found that the target classes are partially balanced. If the target classes were highly unbalanced, then this dataset could not be used because the results of supervised learning algorithms used to make predictions would skew towards the class with the class with higher percentage of records.

```
Accepted coupons: 7210 56.843 %
Rejected coupons: 5474 43.157 %
```

Figure 30: Distribution of Target Classes

It should also be noted that this dataset comes with many missing values. Hence, this dataset requires preprocessing before it can be analyzed with machine learning algorithms. Using the code provided by Niralidedaniya (2023), we can see the features which have missing values.


```

Is there any missing value present or not? True
missing_count missing_percentage
destination      0      0.000000
passanger        0      0.000000
weather          0      0.000000
temperature      0      0.000000
time             0      0.000000
coupon           0      0.000000
expiration       0      0.000000
gender           0      0.000000
age              0      0.000000
maritalStatus    0      0.000000
has_children     0      0.000000
education        0      0.000000
occupation       0      0.000000
income           0      0.000000
car              12576    99.148534
Bar              107      0.843582
CoffeeHouse      217      1.710817
CarryAway        151      1.190476
RestaurantLessThan20 130      1.024913
Restaurant20To50 189      1.490066
toCoupon_GEQ5min 0      0.000000
toCoupon_GEQ15min 0      0.000000
toCoupon_GEQ25min 0      0.000000
direction_same   0      0.000000
direction_opp    0      0.000000
Y                0      0.000000

```

Figure 31: Distribution of missing values in the dataset

There are too many missing values for car feature. To avoid unreliable results, the car feature will be dropped from the dataset. toCoupon_GEQ5min has the same value for all rows. This does not add any useful information to the analysis. toCoupon_GEQ5min will be dropped from

the dataset. Occupation feature has too many categories. This creates too much noise and distracts from clear analysis. Occupation feature will be dropped from the dataset.

```
#Drop column car as it has too many missing values
#Drop column toCoupon_GEQ5min because there is no variability in its value
#Drop column occupation because it has too many categories that leads to a lot of noise
in_vehicle_coupon_data = in_vehicle_coupon_data.drop(['car', 'occupation', 'toCoupon_GEQ5min'], axis=1)
in_vehicle_coupon_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 12684 entries, 0 to 12683
Data columns (total 23 columns):
#   Column                Non-Null Count  Dtype
---  -
0   destination            12684 non-null  object
1   passanger              12684 non-null  object
2   weather                12684 non-null  object
3   temperature            12684 non-null  int64
4   time                   12684 non-null  object
5   coupon                 12684 non-null  object
6   expiration             12684 non-null  object
7   gender                 12684 non-null  object
8   age                   12684 non-null  object
9   maritalStatus          12684 non-null  object
10  has_children           12684 non-null  int64
11  education              12684 non-null  object
12  income                 12684 non-null  object
13  Bar                    12577 non-null  object
14  CoffeeHouse            12467 non-null  object
15  CarryAway              12533 non-null  object
16  RestaurantLessThan20   12554 non-null  object
17  Restaurant20To50       12495 non-null  object
18  toCoupon_GEQ15min      12684 non-null  int64
19  toCoupon_GEQ25min      12684 non-null  int64
20  direction_same         12684 non-null  int64
21  direction_opp          12684 non-null  int64
22  y                      12684 non-null  int64
dtypes: int64(7), object(16)
memory usage: 2.2+ MB
```

Figure 32: Information of dataset after dropping features: car, toCoupon_GEQ5min, occupation

From Table 1, Atiq et al. (2022) demonstrated pretty good accuracy when using frequent value imputation for missing values in the dataset along with the algorithms: Random Forest, Decision Tree, Logistic Regression, Gradient Boosted Tree, Naïve Bayes and Deep Learning. Therefore, this project is going to use frequent value imputation too for the remaining missing values. This can be easily done with the help of code snippet provided by Niralidedaniya (2022).

```
[13] # frequent value / mode imputation for missing values in data. This code snippet has been taken from Miralidedaniya (2022).
in_vehicle_coupon_data['Bar']=in_vehicle_coupon_data['Bar'].fillna(in_vehicle_coupon_data['Bar'].value_counts().index[0])
in_vehicle_coupon_data['CoffeeHouse']=in_vehicle_coupon_data['CoffeeHouse'].fillna(in_vehicle_coupon_data['CoffeeHouse'].value_counts().index[0])
in_vehicle_coupon_data['CarryAway']=in_vehicle_coupon_data['CarryAway'].fillna(in_vehicle_coupon_data['CarryAway'].value_counts().index[0])
in_vehicle_coupon_data['RestaurantLessThan20']=in_vehicle_coupon_data['RestaurantLessThan20'].fillna(in_vehicle_coupon_data['RestaurantLessThan20'].value_counts().index[0])
in_vehicle_coupon_data['Restaurant20To50']=in_vehicle_coupon_data['Restaurant20To50'].fillna(in_vehicle_coupon_data['Restaurant20To50'].value_counts().index[0])
#lets check for missing values again
print('Is there any missing value present?',in_vehicle_coupon_data.isnull().values.any())
```

Is there any missing value present? False

Figure 33: Mode / frequent value imputation for missing values

Using ordinal encoding for the categorical values in the categorical features, the covariance matrix was generated. Strangely, direction_same feature has same values as direction_opp feature in the covariance matrix.

```
from sklearn.preprocessing import OrdinalEncoder
encoder = OrdinalEncoder()
encoded_in_vehicle_coupon_data = encoder.fit_transform(in_vehicle_coupon_data)
encoded_in_vehicle_coupon_data = pd.DataFrame(encoded_in_vehicle_coupon_data, columns=in_vehicle_coupon_data.columns)
encoded_in_vehicle_coupon_data.cov()
```

	destination	passanger	weather	temperature	time	coupon	expiration	gender	age	maritalStatus	...	Bar	CoffeeHouse	CarryAway	Restaurant
destination	0.504658	-0.078836	-0.035245	-0.015041	0.420743	-0.010986	-0.011497	0.002407	-0.003339	0.001797	...	-0.005564	-0.007850	-0.005959	
passanger	-0.078836	0.887317	0.035703	0.040508	-0.543328	0.023840	0.034968	-0.009838	0.010415	-0.016771	...	0.022005	-0.012488	-0.017730	
weather	-0.035245	0.035703	0.401445	0.210926	-0.021266	0.125609	0.005569	-0.008552	-0.027806	-0.008874	...	0.010249	-0.005865	-0.026723	
temperature	-0.015041	0.040508	0.210926	0.587031	-0.065117	0.133050	0.047207	-0.009767	-0.046974	0.002288	...	0.010868	0.008828	-0.027345	
time	0.420743	-0.543328	-0.021266	-0.065117	2.072183	0.092426	-0.060532	-0.002523	-0.025719	0.005786	...	-0.009895	-0.011495	0.001266	
coupon	-0.010986	0.023840	0.125609	0.133050	0.092426	1.818577	0.099355	0.004809	0.006944	-0.000149	...	-0.015862	-0.004406	-0.004440	
expiration	-0.011497	0.034968	0.005569	0.047207	-0.060532	0.099355	0.246532	-0.000314	0.007207	-0.005098	...	-0.006669	-0.009675	-0.002079	
gender	0.002407	-0.009838	-0.008552	-0.009767	-0.002523	0.004809	-0.000314	0.249842	-0.067526	0.023960	...	-0.120201	0.045669	-0.009576	
age	-0.003339	0.010415	-0.027806	-0.046974	-0.025719	0.006944	0.007207	-0.067526	4.950184	-0.328927	...	0.712399	-0.002177	0.216170	
maritalStatus	0.001797	-0.016771	-0.008874	0.002288	0.005786	-0.000149	-0.005098	0.023960	-0.328927	0.693751	...	-0.118302	0.014938	-0.022649	
has_children	-0.002347	0.016028	0.003950	-0.007441	-0.005145	-0.006923	0.003918	-0.039384	0.335703	-0.177913	...	0.127863	-0.010256	0.003538	
education	0.011758	0.001656	0.015185	0.022296	-0.008154	-0.001837	-0.008597	0.014140	0.363653	0.107978	...	0.116209	0.011414	0.049950	
income	-0.025464	-0.005402	-0.051117	-0.047347	-0.022016	0.002894	-0.013671	0.032841	0.267308	0.150888	...	0.053494	-0.134355	0.120053	
Bar	-0.005564	0.022005	0.010249	0.010868	-0.009895	-0.015862	-0.006669	-0.120201	0.712399	-0.118302	...	2.407485	0.378947	0.134812	
CoffeeHouse	-0.007850	-0.012488	-0.005865	0.008828	-0.011495	-0.004406	-0.009675	0.045669	-0.002177	0.014938	...	0.378947	2.365560	0.195541	
CarryAway	-0.005959	-0.017730	-0.026723	-0.027345	0.001266	-0.004440	-0.002079	-0.009576	0.216170	-0.022649	...	0.134812	0.195541	1.194358	
RestaurantLessThan20	0.002886	-0.033238	-0.003905	-0.001147	-0.004725	0.016685	-0.006039	0.025848	-0.074016	-0.001212	...	0.139119	0.306658	0.113426	
Restaurant20To50	-0.000269	-0.042455	0.005444	0.002082	0.019281	0.011179	-0.000336	-0.000155	-0.004063	0.055322	...	0.375264	0.321059	0.193181	
toCoupon_GEQ15min	0.049593	0.030170	-0.038262	-0.059057	0.005156	-0.088045	0.010531	-0.001743	0.029335	-0.020447	...	0.015285	-0.003662	0.002056	
toCoupon_GEQ25min	0.045391	-0.060297	-0.041579	-0.053675	0.136276	-0.049269	-0.005304	0.000444	-0.000045	0.001348	...	0.002514	0.001584	0.001954	
direction_same	-0.024310	-0.103995	0.004609	0.030548	0.184128	-0.040432	0.006848	-0.000923	-0.007565	0.005645	...	-0.000763	0.011396	0.003997	
direction_opp	0.024310	0.103995	-0.004609	-0.030548	-0.184128	0.040432	-0.006848	0.000923	0.007565	-0.005645	...	0.000763	-0.011396	-0.003997	
Y	-0.000671	0.024082	0.031006	0.023241	-0.033780	0.064804	-0.031952	0.010886	-0.038836	0.010348	...	-0.058434	-0.110180	-0.026371	

Figure 34: Covariance matrix on the dataset after doing ordinal encoding on the categorical values of the dataset.

Since `direction_same` has same covariance values as `direction_opp`. It makes sense to just have one on them and reduce the noise.

```
in_vehicle_coupon_data = in_vehicle_coupon_data.drop(['direction_opp'], axis=1)
in_vehicle_coupon_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 12684 entries, 0 to 12683
Data columns (total 22 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   destination                          12684 non-null  object
1   passenger                            12684 non-null  object
2   weather                             12684 non-null  object
3   temperature                         12684 non-null  int64
4   time                                12684 non-null  object
5   coupon                             12684 non-null  object
6   expiration                          12684 non-null  object
7   gender                              12684 non-null  object
8   age                                 12684 non-null  object
9   maritalStatus                      12684 non-null  object
10  has_children                        12684 non-null  int64
11  education                          12684 non-null  object
12  income                             12684 non-null  object
13  Bar                                 12684 non-null  object
14  CoffeeHouse                        12684 non-null  object
15  CarryAway                          12684 non-null  object
16  RestaurantLessThan20               12684 non-null  object
17  Restaurant20To50                   12684 non-null  object
18  toCoupon_GEQ15min                  12684 non-null  int64
19  toCoupon_GEQ25min                  12684 non-null  int64
20  direction_same                     12684 non-null  int64
21  Y                                   12684 non-null  int64
dtypes: int64(6), object(16)
memory usage: 2.1+ MB
```

Figure 35: Dataset information after data cleaning

With all the steps above, the data is cleaned and ready to use. We can also generate a correlation matrix.

```
#Get the correlation matrix using ordinal encoding for the categorical values
encoder = OrdinalEncoder()
encoded_in_vehicle_coupon_data = encoder.fit_transform(in_vehicle_coupon_data)
encoded_in_vehicle_coupon_data = pd.DataFrame(encoded_in_vehicle_coupon_data, columns=in_vehicle_coupon_data.columns)
correlation_matrix = encoded_in_vehicle_coupon_data.corr()
print(correlation_matrix)
```

	destination	passanger	weather	temperature	time \
destination	1.000000	-0.117811	-0.078305	-0.027633	0.411437
passanger	-0.117811	1.000000	0.059821	0.056127	-0.400690
weather	-0.078305	0.059821	1.000000	0.434497	-0.023316
temperature	-0.027633	0.056127	0.434497	1.000000	-0.059041
time	0.411437	-0.400690	-0.023316	-0.059041	1.000000
coupon	-0.011468	0.018767	0.147008	0.128771	0.047612
expiration	-0.032594	0.074764	0.017702	0.124090	-0.084691
gender	0.006779	-0.020896	-0.027003	-0.025504	-0.003507
age	-0.002112	0.004969	-0.019725	-0.027556	-0.008030
maritalStatus	0.003036	-0.021376	-0.016816	0.003585	0.004826
has_children	-0.006707	0.034542	0.012657	-0.019716	-0.007256
education	0.008793	0.000934	0.012733	0.015460	-0.003009
income	-0.014554	-0.002329	-0.032758	-0.025091	-0.006210
Bar	-0.005048	0.015055	0.010426	0.009142	-0.004430
CoffeeHouse	-0.007185	-0.008619	-0.006018	0.007491	-0.005192
CarryAway	-0.007676	-0.017223	-0.038593	-0.032657	0.000804
RestaurantLessThan20	0.003497	-0.030378	-0.005306	-0.001289	-0.002826
Restaurant20To50	-0.000255	-0.030401	0.005795	0.001833	0.009035
toCoupon_GEQ15min	0.140684	0.064544	-0.121698	-0.155332	0.007218
toCoupon_GEQ25min	0.197240	-0.197595	-0.202572	-0.216254	0.292231
direction_same	-0.083328	-0.268830	0.017712	0.097085	0.311467
Y	-0.001906	0.051614	0.098800	0.061240	-0.047377

	coupon	expiration	gender	age	maritalStatus \
destination	-0.011468	-0.032594	0.006779	-0.002112	0.003036
passanger	0.018767	0.074764	-0.020896	0.004969	-0.021376
weather	0.147008	0.017702	-0.027003	-0.019725	-0.016816
temperature	0.128771	0.124090	-0.025504	-0.027556	0.003585
time	0.047612	-0.084691	-0.003507	-0.008030	0.004826
coupon	1.000000	0.148383	0.007134	0.002314	-0.000132
expiration	0.148383	1.000000	-0.001264	0.006523	-0.012328
gender	0.007134	-0.001264	1.000000	-0.060720	0.057552
age	0.002314	0.006523	-0.060720	1.000000	-0.177495
maritalStatus	-0.000132	-0.012328	0.057552	-0.177495	1.000000
has_children	-0.010422	0.016020	-0.159956	0.306306	-0.433628
education	-0.000724	-0.009198	0.015029	0.086833	0.068872
income	0.000871	-0.011180	0.026677	0.048782	0.073555
Bar	-0.007581	-0.008656	-0.154986	0.206363	-0.091539
CoffeeHouse	-0.002124	-0.012669	0.059404	-0.000636	0.011661
CarryAway	-0.003012	-0.003832	-0.017530	0.088903	-0.024881
RestaurantLessThan20	0.010651	-0.010471	0.044519	-0.028640	-0.001253
Restaurant20To50	0.005592	-0.000457	-0.000209	-0.001232	0.044802
toCoupon_GEQ15min	-0.131571	0.042740	-0.007028	0.026571	-0.049471
toCoupon_GEQ25min	-0.112780	-0.032977	0.002743	-0.000063	0.004997
direction_same	-0.073007	0.033584	-0.004496	-0.008279	0.016504
Y	0.097019	-0.129920	0.043969	-0.035241	0.025083

	...	income	Bar	CoffeeHouse	CarryAway	\
destination	...	-0.014554	-0.005048	-0.007185	-0.007676	
passanger	...	-0.002329	0.015055	-0.008619	-0.017223	
weather	...	-0.032758	0.010426	-0.006018	-0.038593	
temperature	...	-0.025091	0.009142	0.007491	-0.032657	
time	...	-0.006210	-0.004430	-0.005192	0.000804	
coupon	...	0.000871	-0.007581	-0.002124	-0.003012	
expiration	...	-0.011180	-0.008656	-0.012669	-0.003832	
gender	...	0.026677	-0.154986	0.059404	-0.017530	
age	...	0.048782	0.206363	-0.000636	0.088903	
maritalStatus	...	0.073555	-0.091539	0.011661	-0.024881	
has_children	...	-0.000569	0.167292	-0.013537	0.006573	
education	...	0.001157	0.039789	0.003942	0.024282	
income	...	1.000000	0.013998	-0.035469	0.044603	
Bar	...	0.013998	1.000000	0.158793	0.079502	
CoffeeHouse	...	-0.035469	0.158793	1.000000	0.116333	
CarryAway	...	0.044603	0.079502	0.116333	1.000000	
RestaurantLessThan20	...	0.053753	0.077191	0.171650	0.089352	
Restaurant20To50	...	0.018483	0.163138	0.140805	0.119233	
toCoupon_GEQ15min	...	0.006196	0.019852	-0.004799	0.003791	
toCoupon_GEQ25min	...	0.003940	0.005002	0.003179	0.005519	
direction_same	...	0.017048	-0.001197	0.018042	0.008906	
Y	...	-0.023949	-0.076033	-0.144629	-0.048717	

	RestaurantLessThan20	Restaurant20To50	\
destination	0.003497	-0.000255	
passanger	-0.030378	-0.030401	
weather	-0.005306	0.005795	
temperature	-0.001289	0.001833	
time	-0.002826	0.009035	
coupon	0.010651	0.005592	
expiration	-0.010471	-0.000457	
gender	0.044519	-0.000209	
age	-0.028640	-0.001232	
maritalStatus	-0.001253	0.044802	
has_children	-0.074252	-0.000614	
education	0.026678	0.081681	
income	0.053753	0.018483	
Bar	0.077191	0.163138	
CoffeeHouse	0.171650	0.140805	
CarryAway	0.089352	0.119233	
RestaurantLessThan20	1.000000	0.111437	
Restaurant20To50	0.111437	1.000000	
toCoupon_GEQ15min	0.004797	-0.001461	
toCoupon_GEQ25min	-0.003029	0.008245	
direction_same	0.003103	0.007554	
Y	-0.011137	-0.056268	

	toCoupon_GEQ15min	toCoupon_GEQ25min	direction_same \
destination	0.140684	0.197240	-0.083328
passanger	0.064544	-0.197595	-0.268830
weather	-0.121698	-0.202572	0.017712
temperature	-0.155332	-0.216254	0.097085
time	0.007218	0.292231	0.311467
coupon	-0.131571	-0.112780	-0.073007
expiration	0.042740	-0.032977	0.033584
gender	-0.007028	0.002743	-0.004496
age	0.026571	-0.000063	-0.008279
maritalStatus	-0.049471	0.004997	0.016504
has_children	0.078211	-0.013722	-0.031620
education	-0.019017	-0.010812	0.001205
income	0.006196	0.003940	0.017048
Bar	0.019852	0.005002	-0.001197
CoffeeHouse	-0.004799	0.003179	0.018042
CarryAway	0.003791	0.005519	0.008906
RestaurantLessThan20	0.004797	-0.003029	0.003103
Restaurant20To50	-0.001461	0.008245	0.007554
toCoupon_GEQ15min	1.000000	0.324984	-0.303533
toCoupon_GEQ25min	0.324984	1.000000	-0.192319
direction_same	-0.303533	-0.192319	1.000000
Y	-0.081602	-0.103633	0.014570

	Y
destination	-0.001906
passanger	0.051614
weather	0.098800
temperature	0.061240
time	-0.047377
coupon	0.097019
expiration	-0.129920
gender	0.043969
age	-0.035241
maritalStatus	0.025083
has_children	-0.045557
education	0.043023
income	-0.023949
Bar	-0.076033
CoffeeHouse	-0.144629
CarryAway	-0.048717
RestaurantLessThan20	-0.011137
Restaurant20To50	-0.056268
toCoupon_GEQ15min	-0.081602
toCoupon_GEQ25min	-0.103633
direction_same	0.014570
Y	1.000000

Figure 36: Correlation matrix after cleaning In Vehicle Coupon dataset

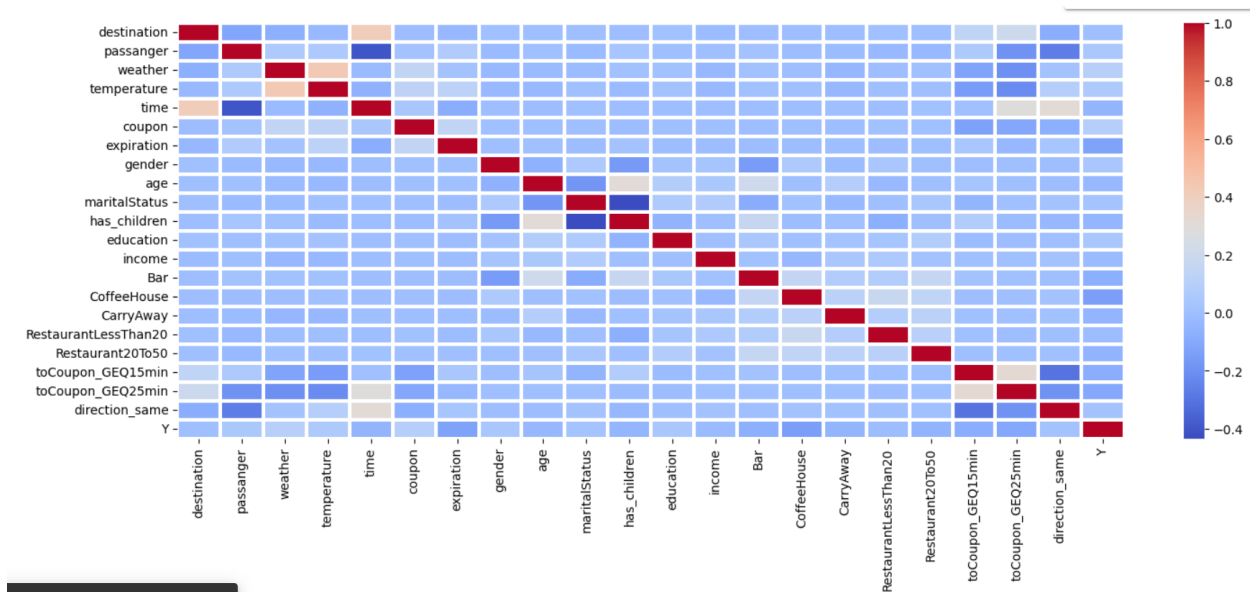


Figure 37: Heat map of correlation matrix

A heatmap to visualize the correlation between features is also effective.

There seems to be a correlation between time and destination, between temperature and weather, between marital status and has children, between passenger and time, between to coupon GEQ 15 min and same direction. **There seems to be correlation of the following features with the target Y which determines whether the coupon was accepted or rejected: expiration, CoffeeHouse, toCoupon GEQ15min, toCoupon GEQ25min. In the heat map, the darker the color, the stronger the correlation. This also corresponds with the higher values in the correlation matrix.**

Dimensionality Reduction

Dimensionality reduction was done using 3 methods: Stepwise Regression, Forward Feature Selection and Backward Feature Elimination methods learnt by us in *CMTH 642 – Data Analytics: Advanced Methods*. This will help eliminate some noise. **Dimensionality reduction is the major contribution of this project to this dataset.** Past work as seen in the Literature Review did not do dimensionality reduction on this dataset.

Stepwise Regression

Alpha of 0.05 was chosen. Stepwise Regression for dimensionality reduction was done using Python. Coefficients having a p-value of 0.05 or less will be statistically significant.

Iteration 1: Bar has highest p-value of 0.779. Bar will be dropped. It is the least statistically significant.

OLS Regression Results						
=====						
Dep. Variable:	Y	R-squared (uncentered):	0.587			
Model:	OLS	Adj. R-squared (uncentered):	0.586			
Method:	Least Squares	F-statistic:	856.7			
Date:	Thu, 28 Nov 2024	Prob (F-statistic):	0.00			
Time:	01:41:07	Log-Likelihood:	-8808.8			
No. Observations:	12684	AIC:	1.766e+04			
Df Residuals:	12663	BIC:	1.782e+04			
Df Model:	21					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

destination	0.0569	0.007	8.321	0.000	0.044	0.070
passanger	0.0471	0.005	9.367	0.000	0.037	0.057
weather	0.1170	0.007	16.382	0.000	0.103	0.131
temperature	0.0266	0.006	4.160	0.000	0.014	0.039
time	-0.0096	0.004	-2.434	0.015	-0.017	-0.002
coupon	0.0541	0.003	16.362	0.000	0.048	0.061
expiration	-0.1407	0.009	-15.771	0.000	-0.158	-0.123
gender	0.0866	0.009	10.012	0.000	0.070	0.104
age	0.0025	0.002	1.189	0.234	-0.002	0.007
maritalStatus	0.0635	0.005	12.183	0.000	0.053	0.074
has_children	0.0444	0.010	4.483	0.000	0.025	0.064
education	0.0200	0.002	8.762	0.000	0.016	0.025
income	0.0014	0.002	0.820	0.412	-0.002	0.005
Bar	0.0008	0.003	0.281	0.779	-0.005	0.007
CoffeeHouse	-0.0363	0.003	-12.574	0.000	-0.042	-0.031
CarryAway	0.0014	0.004	0.354	0.723	-0.006	0.009
RestaurantLessThan20	0.0150	0.004	3.935	0.000	0.008	0.022
Restaurant20To50	-0.0019	0.003	-0.638	0.524	-0.008	0.004
toCoupon_GEQ15min	0.0237	0.009	2.502	0.012	0.005	0.042
toCoupon_GEQ25min	-0.0321	0.016	-2.039	0.041	-0.063	-0.001
direction_same	0.1026	0.013	8.132	0.000	0.078	0.127
=====						
Omnibus:	77704.989	Durbin-Watson:	1.682			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	1296.009			
Skew:	-0.228	Prob(JB):	3.76e-282			
Kurtosis:	1.502	Cond. No.	31.5			
=====						

Figure 38: First iteration of stepwise regression

Iteration 2: CarryAway is least statistically significant because it has highest p-value of 0.713 which is greater than alpha of 0.05. CarryAway should be dropped.

OLS Regression Results						
=====						
Dep. Variable:	Y	R-squared (uncentered):	0.587			
Model:	OLS	Adj. R-squared (uncentered):	0.586			
Method:	Least Squares	F-statistic:	899.5			
Date:	Thu, 28 Nov 2024	Prob (F-statistic):	0.00			
Time:	01:41:07	Log-Likelihood:	-8808.9			
No. Observations:	12684	AIC:	1.766e+04			
Df Residuals:	12664	BIC:	1.781e+04			
Df Model:	20					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

destination	0.0570	0.007	8.332	0.000	0.044	0.070
passanger	0.0472	0.005	9.400	0.000	0.037	0.057
weather	0.1172	0.007	16.463	0.000	0.103	0.131
temperature	0.0266	0.006	4.172	0.000	0.014	0.039
time	-0.0095	0.004	-2.426	0.015	-0.017	-0.002
coupon	0.0541	0.003	16.380	0.000	0.048	0.061
expiration	-0.1407	0.009	-15.770	0.000	-0.158	-0.123
gender	0.0863	0.009	10.041	0.000	0.069	0.103
age	0.0026	0.002	1.265	0.206	-0.001	0.007
maritalStatus	0.0636	0.005	12.229	0.000	0.053	0.074
has_children	0.0448	0.010	4.567	0.000	0.026	0.064
education	0.0201	0.002	8.790	0.000	0.016	0.025
income	0.0014	0.002	0.833	0.405	-0.002	0.005
CoffeeHouse	-0.0362	0.003	-12.713	0.000	-0.042	-0.031
CarryAway	0.0015	0.004	0.368	0.713	-0.006	0.009
RestaurantLessThan20	0.0150	0.004	3.966	0.000	0.008	0.022
Restaurant20To50	-0.0018	0.003	-0.599	0.549	-0.008	0.004
toCoupon_GEQ15min	0.0238	0.009	2.519	0.012	0.005	0.042
toCoupon_GEQ25min	-0.0319	0.016	-2.032	0.042	-0.063	-0.001
direction_same	0.1027	0.013	8.147	0.000	0.078	0.127
=====						
Omnibus:	77832.632	Durbin-Watson:	1.682			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	1295.135			
Skew:	-0.228	Prob(JB):	5.82e-282			
Kurtosis:	1.503	Cond. No.	29.7			
=====						

Figure 39: Second iteration of stepwise regression

Iteration 3: Restaurant20To50 is least statistically significant because it has highest p-value of 0.576 which is greater than alpha of 0.05. Restaurant20To50 should be dropped.

OLS Regression Results						
=====						
Dep. Variable:	Y	R-squared (uncentered):	0.587			
Model:	OLS	Adj. R-squared (uncentered):	0.586			
Method:	Least Squares	F-statistic:	946.9			
Date:	Thu, 28 Nov 2024	Prob (F-statistic):	0.00			
Time:	01:41:07	Log-Likelihood:	-8808.9			
No. Observations:	12684	AIC:	1.766e+04			
Df Residuals:	12665	BIC:	1.780e+04			
Df Model:	19					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

destination	0.0570	0.007	8.337	0.000	0.044	0.070
passanger	0.0472	0.005	9.409	0.000	0.037	0.057
weather	0.1173	0.007	16.477	0.000	0.103	0.131
temperature	0.0266	0.006	4.171	0.000	0.014	0.039
time	-0.0095	0.004	-2.420	0.016	-0.017	-0.002
coupon	0.0541	0.003	16.396	0.000	0.048	0.061
expiration	-0.1406	0.009	-15.767	0.000	-0.158	-0.123
gender	0.0863	0.009	10.040	0.000	0.069	0.103
age	0.0027	0.002	1.312	0.189	-0.001	0.007
maritalStatus	0.0637	0.005	12.246	0.000	0.053	0.074
has_children	0.0448	0.010	4.570	0.000	0.026	0.064
education	0.0201	0.002	8.805	0.000	0.016	0.025
income	0.0015	0.002	0.859	0.390	-0.002	0.005
CoffeeHouse	-0.0361	0.003	-12.757	0.000	-0.042	-0.031
RestaurantLessThan20	0.0151	0.004	4.005	0.000	0.008	0.023
Restaurant20To50	-0.0016	0.003	-0.559	0.576	-0.007	0.004
toCoupon_GEQ15min	0.0240	0.009	2.532	0.011	0.005	0.042
toCoupon_GEQ25min	-0.0319	0.016	-2.030	0.042	-0.063	-0.001
direction_same	0.1029	0.013	8.158	0.000	0.078	0.128
=====						
Omnibus:	77965.192	Durbin-Watson:	1.682			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	1294.195			
Skew:	-0.228	Prob(JB):	9.31e-282			
Kurtosis:	1.503	Cond. No.	29.4			
=====						

Figure 40: Third iteration of stepwise regression

Iteration 4: Income is least statistically significant because it has highest p-value of 0.405 which is greater than alpha of 0.05. Income should be dropped.

OLS Regression Results						
Dep. Variable:	Y	R-squared (uncentered):	0.587			
Model:	OLS	Adj. R-squared (uncentered):	0.586			
Method:	Least Squares	F-statistic:	999.6			
Date:	Thu, 28 Nov 2024	Prob (F-statistic):	0.00			
Time:	01:41:07	Log-Likelihood:	-8809.1			
No. Observations:	12684	AIC:	1.765e+04			
Df Residuals:	12666	BIC:	1.779e+04			
Df Model:	18					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
destination	0.0569	0.007	8.326	0.000	0.044	0.070
passanger	0.0472	0.005	9.402	0.000	0.037	0.057
weather	0.1169	0.007	16.487	0.000	0.103	0.131
temperature	0.0265	0.006	4.161	0.000	0.014	0.039
time	-0.0096	0.004	-2.436	0.015	-0.017	-0.002
coupon	0.0541	0.003	16.388	0.000	0.048	0.061
expiration	-0.1407	0.009	-15.782	0.000	-0.158	-0.123
gender	0.0862	0.009	10.031	0.000	0.069	0.103
age	0.0026	0.002	1.300	0.193	-0.001	0.007
maritalStatus	0.0633	0.005	12.288	0.000	0.053	0.073
has_children	0.0444	0.010	4.539	0.000	0.025	0.064
education	0.0200	0.002	8.798	0.000	0.016	0.024
income	0.0014	0.002	0.832	0.405	-0.002	0.005
CoffeeHouse	-0.0363	0.003	-13.018	0.000	-0.042	-0.031
RestaurantLessThan20	0.0149	0.004	3.967	0.000	0.008	0.022
toCoupon_GEQ15min	0.0237	0.009	2.510	0.012	0.005	0.042
toCoupon_GEQ25min	-0.0321	0.016	-2.042	0.041	-0.063	-0.001
direction_same	0.1027	0.013	8.146	0.000	0.078	0.127
Omnibus:	77733.204	Durbin-Watson:	1.682			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	1296.124			
Skew:	-0.229	Prob(JB):	3.55e-282			
Kurtosis:	1.502	Cond. No.	28.1			

Figure 41: Fourth iteration of stepwise regression

Iteration 5: Age is least statistically significant because it has highest p-value of 0.166 which is greater than alpha of 0.05. age should be dropped.

OLS Regression Results						
=====						
Dep. Variable:	Y	R-squared (uncentered):		0.587		
Model:	OLS	Adj. R-squared (uncentered):		0.586		
Method:	Least Squares	F-statistic:		1058.		
Date:	Thu, 28 Nov 2024	Prob (F-statistic):		0.00		
Time:	01:41:07	Log-Likelihood:		-8809.4		
No. Observations:	12684	AIC:		1.765e+04		
Df Residuals:	12667	BIC:		1.778e+04		
Df Model:	17					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

destination	0.0570	0.007	8.345	0.000	0.044	0.070
passanger	0.0474	0.005	9.456	0.000	0.038	0.057
weather	0.1173	0.007	16.564	0.000	0.103	0.131
temperature	0.0266	0.006	4.172	0.000	0.014	0.039
time	-0.0095	0.004	-2.419	0.016	-0.017	-0.002
coupon	0.0542	0.003	16.455	0.000	0.048	0.061
expiration	-0.1407	0.009	-15.779	0.000	-0.158	-0.123
gender	0.0867	0.009	10.112	0.000	0.070	0.103
age	0.0028	0.002	1.384	0.166	-0.001	0.007
maritalStatus	0.0641	0.005	12.716	0.000	0.054	0.074
has_children	0.0451	0.010	4.632	0.000	0.026	0.064
education	0.0200	0.002	8.825	0.000	0.016	0.024
CoffeeHouse	-0.0364	0.003	-13.027	0.000	-0.042	-0.031
RestaurantLessThan20	0.0152	0.004	4.061	0.000	0.008	0.023
toCoupon_GEQ15min	0.0242	0.009	2.565	0.010	0.006	0.043
toCoupon_GEQ25min	-0.0318	0.016	-2.023	0.043	-0.063	-0.001
direction_same	0.1033	0.013	8.208	0.000	0.079	0.128
=====						
Omnibus:	78102.225	Durbin-Watson:		1.682		
Prob(Omnibus):	0.000	Jarque-Bera (JB):		1293.019		
Skew:	-0.228	Prob(JB):		1.68e-281		
Kurtosis:	1.504	Cond. No.		24.7		
=====						

Figure 42: Fifth iteration of stepwise regression

After 5 iterations, all coefficients are statistically significant.

OLS Regression Results						
=====						
Dep. Variable:	Y	R-squared (uncentered):	0.587			
Model:	OLS	Adj. R-squared (uncentered):	0.586			
Method:	Least Squares	F-statistic:	1124.			
Date:	Thu, 28 Nov 2024	Prob (F-statistic):	0.00			
Time:	01:41:08	Log-Likelihood:	-8810.4			
No. Observations:	12684	AIC:	1.765e+04			
Df Residuals:	12668	BIC:	1.777e+04			
Df Model:	16					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

destination	0.0573	0.007	8.395	0.000	0.044	0.071
passanger	0.0476	0.005	9.511	0.000	0.038	0.057
weather	0.1179	0.007	16.677	0.000	0.104	0.132
temperature	0.0266	0.006	4.178	0.000	0.014	0.039
time	-0.0094	0.004	-2.393	0.017	-0.017	-0.002
coupon	0.0545	0.003	16.551	0.000	0.048	0.061
expiration	-0.1405	0.009	-15.755	0.000	-0.158	-0.123
gender	0.0869	0.009	10.147	0.000	0.070	0.104
maritalStatus	0.0644	0.005	12.772	0.000	0.054	0.074
has_children	0.0496	0.009	5.403	0.000	0.032	0.068
education	0.0205	0.002	9.148	0.000	0.016	0.025
CoffeeHouse	-0.0362	0.003	-12.978	0.000	-0.042	-0.031
RestaurantLessThan20	0.0153	0.004	4.082	0.000	0.008	0.023
toCoupon_GEQ15min	0.0248	0.009	2.636	0.008	0.006	0.043
toCoupon_GEQ25min	-0.0313	0.016	-1.995	0.046	-0.062	-0.001
direction_same	0.1040	0.013	8.269	0.000	0.079	0.129
=====						
Omnibus:	78704.924	Durbin-Watson:	1.682			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	1288.500			
Skew:	-0.228	Prob(JB):	1.61e-280			
Kurtosis:	1.507	Cond. No.	22.4			
=====						

Figure 43: All coefficients are statistically significant

Therefore, age, income, Restaurant20To50, CarryAway and Bar features are dropped to create a dataset that is dimensionally reduced due to stepwise regression. This leads to a dataset with 16 features.

Forward feature selection

Python was used to do forward feature selection. 16 features were selected for forward feature selection because stepwise regression resulted with a dataset with 16 features. This was done for consistency. The focus was on reducing the negative mean squared error. Forward feature selection resulted in a dataset where maritalStatus, Bar, RestaurantLessThan20, temperature and direction_same features were dropped.

Stepwise Regression and Forward Feature selection for dimensionality reduction gave different results resulting in different columns being dropped.

Backward feature elimination

Python was used to do backward feature elimination. 16 features were selected for backward feature elimination because stepwise regression resulted with a dataset with 16 features. This was done for consistency. The focus was on reducing the negative mean squared error. Backward feature elimination resulted in a dataset where maritalStatus, Bar, RestaurantLessThan20, temperature and direction_same features were dropped.

It has hence observed that forward feature selection and backward feature elimination resulted in the same datasets. Thus, there are 2 dimensionally reduced datasets to work with: the dataset reduced by stepwise regression and the dataset reduced by forward feature selection/backward feature elimination. Dimensionality reduction is the major contribution of towards this dataset compared to past work done on this dataset.

Classification Algorithms and Cross Validation

The performance of 5 classification algorithms was evaluated over each of the 2 dimensionally reduced datasets:

1. Random Forest
2. Logistic Regression
3. k Nearest Neighbours (k-NN)
4. Naïve Bayes
5. Decision Tree

Cross Validation was used and the dataset was split into test and training sets. 20% of dataset was set for testing each time for consistency. The classification algorithms were trained over the training part of the datasets and then evaluated over the testing part of the datasets. Accuracy, classification report and Area under Curve (AUC) were gathered each time.

Result of Random Forest on stepwise regression reduced dataset:

Accuracy: 0.6862435947970044

Classification Report:

	precision	recall	f1-score	support
0.0	0.64	0.61	0.63	1087
1.0	0.72	0.74	0.73	1450
accuracy			0.69	2537
macro avg	0.68	0.68	0.68	2537
weighted avg	0.68	0.69	0.69	2537

AUC: 0.7351197538305365

Figure 44: Result of Random Forest on stepwise regression reduced dataset

Result of Random Forest on forward feature selection/backward feature elimination reduced dataset:

Accuracy: 0.7335435553803705

Classification Report:

	precision	recall	f1-score	support
0.0	0.70	0.65	0.68	1087
1.0	0.75	0.79	0.77	1450
accuracy			0.73	2537
macro avg	0.73	0.72	0.73	2537
weighted avg	0.73	0.73	0.73	2537

AUC: 0.7942876629762395

Figure 45: Result of Random Forest on forward feature selection/backward feature elimination reduced dataset

Result of Logistic Regression on stepwise regression reduced dataset:

```
Accuracy: 0.6243594797004336
Classification Report:
              precision    recall  f1-score   support

    0.0         0.59      0.42      0.49       1087
    1.0         0.64      0.77      0.70       1450

 accuracy          0.62       2537
 macro avg         0.61      0.60      0.60       2537
weighted avg         0.62      0.62      0.61       2537

AUC: 0.66001966817879
```

Figure 46: Result of Logistic Regression on stepwise regression reduced dataset

Result of Logistic Regression on forward feature selection/backward feature elimination reduced dataset:

```
Accuracy: 0.6267244777296019
Classification Report:
              precision    recall  f1-score   support

    0.0         0.59      0.42      0.49       1087
    1.0         0.64      0.78      0.70       1450

 accuracy          0.63       2537
 macro avg         0.62      0.60      0.60       2537
weighted avg         0.62      0.63      0.61       2537

AUC: 0.6605948037940551
```

Figure 47: Result of Logistic Regression on forward feature selection/backward feature elimination reduced dataset

Result of k Nearest Neighbours (k-NN) on stepwise regression reduced dataset:

```

Accuracy: 0.6846669294442255
Classification Report:
              precision    recall  f1-score   support

     0.0         0.64      0.59      0.62       1087
     1.0         0.71      0.75      0.73       1450

 accuracy          0.68          0.68          0.68       2537
  macro avg         0.68          0.67          0.67       2537
 weighted avg         0.68          0.68          0.68       2537

AUC: 0.7027078640992291

```

Figure 48: Result of k Nearest Neighbours (k-NN) on stepwise regression reduced dataset

Result of k Nearest Neighbours (k-NN) on forward feature selection/backward feature elimination reduced dataset:

```

Accuracy: 0.6460386283011431
Classification Report:
              precision    recall  f1-score   support

     0.0         0.60      0.54      0.57       1087
     1.0         0.68      0.73      0.70       1450

 accuracy          0.65       2537
  macro avg         0.64       0.63       0.63       2537
 weighted avg         0.64       0.65       0.64       2537

AUC: 0.6774247374932589

```

Figure 49: Result of k Nearest Neighbours (k-NN) on forward feature selection/backward feature elimination reduced dataset

Result of Naïve Bayes on stepwise regression reduced dataset:

```

Accuracy: 0.5880961765865195
Classification Report:
              precision    recall  f1-score   support

     0.0         0.53      0.34      0.41       1087
     1.0         0.61      0.78      0.68       1450

 accuracy          0.59       2537
  macro avg         0.57       0.56       0.55       2537
 weighted avg         0.58       0.59       0.57       2537

AUC: 0.6340354661675602

```

Figure 50: Result of Naïve Bayes on stepwise regression reduced dataset

Result of Naïve Bayes on forward feature selection/backward feature elimination reduced dataset:

```
Accuracy: 0.5904611746156878
Classification Report:
              precision    recall  f1-score   support

     0.0         0.53      0.34      0.42       1087
     1.0         0.61      0.78      0.68       1450

 accuracy          0.59          2537
 macro avg         0.57          2537
 weighted avg      0.58          2537

AUC: 0.6393208133743615
```

Figure 51: Result of Naïve Bayes on forward feature selection/backward feature elimination reduced dataset

Result of Decision Tree on stepwise regression reduced dataset:

```
Accuracy: 0.6582577847851794
Classification Report:
              precision    recall  f1-score   support

     0.0         0.60      0.63      0.61       1087
     1.0         0.71      0.68      0.69       1450

 accuracy          0.66          2537
 macro avg         0.65          2537
 weighted avg      0.66          2537

AUC: 0.6603676680518986
```

Figure 52: Result of Decision Tree on stepwise regression reduced dataset

Result of Decision Tree on forward feature selection/backward feature elimination reduced dataset:

```
Accuracy: 0.6645644461962948
Classification Report:
              precision    recall  f1-score   support

    0.0         0.60      0.64      0.62       1087
    1.0         0.72      0.68      0.70       1450

 accuracy          0.66          2537
 macro avg         0.66          2537
 weighted avg      0.67          2537

AUC: 0.6618358024299718
```

Figure 53: Result of Decision Tree on forward feature selection/backward feature elimination reduced dataset

Comparison and Analysis of results

Accuracy tells what percentage of the predictions are correct (Ofir Shalev (@ofirdi), 2021).

“Precision and Recall are often in tension. That is, improving Precision typically reduces Recall and vice versa” (Ofir Shalev (@ofirdi), 2021). “F1 score combines Recall and Precision to one performance metric. F1 Score is the weighted average of Precision and Recall. Therefore, this score takes both false positives and false negatives into account” (Ofir Shalev (@ofirdi), 2021).

The Receiver Operating Characteristics (ROC) curve is created by plotting the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold settings (Ofir Shalev (@ofirdi), 2021). AUC (Area Under Curve): The model’s performance is determined by looking at the area under the ROC curve (or AUC) (Ofir Shalev (@ofirdi), 2021).

Accuracy and AUC was compared for the 5 classification algorithms run over the 2 dimensionally reduced dataset.

Classification Algorithm	Dimensionality Reduction Method	Accuracy	AUC
Random Forest	Stepwise regression	0.6862	0.7351
	forward feature selection/backward feature elimination	0.7335	0.7943
Logistic Regression	Stepwise regression	0.6244	0.6600
	forward feature selection/backward feature elimination	0.6267	0.6606
k Nearest Neighbours (k-NN)	Stepwise regression	0.6847	0.7027
	forward feature selection/backward feature elimination	0.6460	0.6774
Naïve Bayes	Stepwise regression	0.5881	0.6340
	forward feature selection/backward feature elimination	0.5905	0.6393
Decision Tree	Stepwise regression	0.6602	0.6612
	forward feature selection/backward feature elimination	0.6618	0.6584

Table 3: Accuracy and AUC on Dimensionally reduced data using different classification algorithms

Random forest on dataset reduced using forward feature selection/backward feature elimination had the highest accuracy and AUC value closest to 1. The Area Under Curve (AUC) close to 1, shows the high predictive power. Whereas Naïve Bayes on Stepwise Regression reduced dataset had lowest accuracy and AUC value closest to 0.5. An AUC value closer to 0.5 shows that it is as good as random chance.

Methodology

The methodology being followed is based on the video by Babaoglu. (2018, January 6).

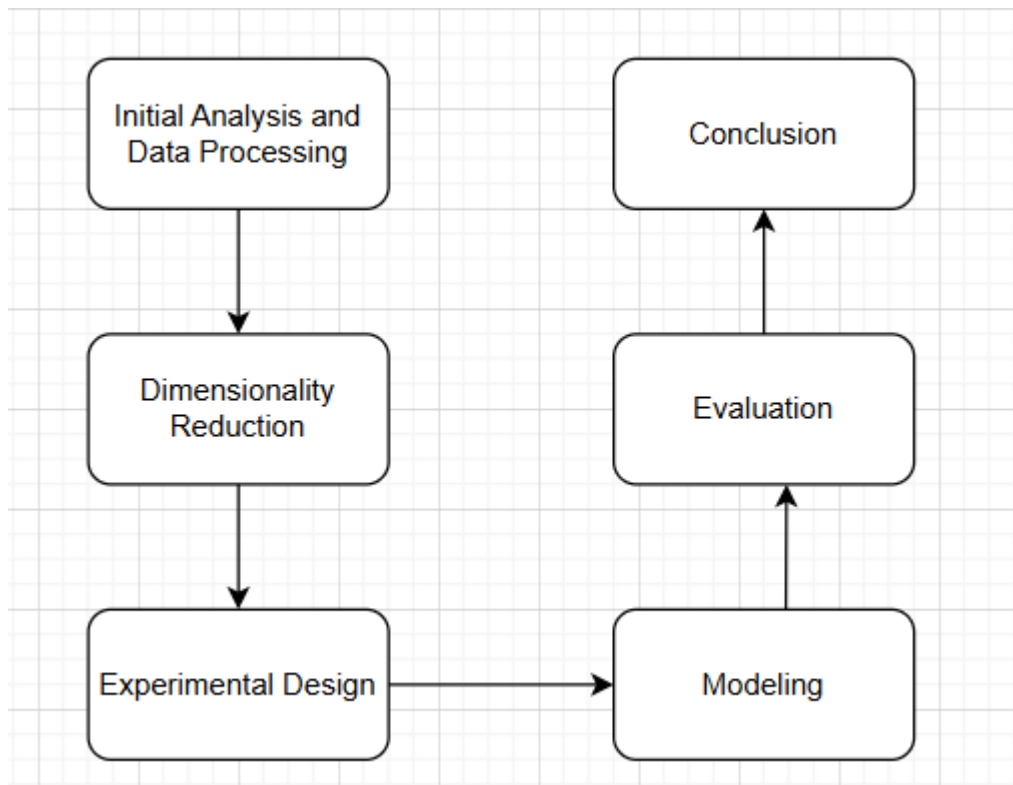


Figure 54: Proposed methodology for this project

In the Initial Analysis, the data will be described and examined. Univariate Analysis will be done to look at the distribution of each attribute. Check for missing values will be conducted. Data Processing will also be done and data will be cleaned.

During dimensionality reduction, a dataset with fewer dimensions using these 3 methods:

Stepwise Regression, Forward Feature Selection and Backward Feature Elimination will be acquired.

Experimental Design will involve splitting data into a training set and test set to run classification algorithms over.

Modeling will involve running Random Forest, Decision Tree, Logistic Regression, Naïve Bayes, k-Nearest Neighbours (k-NN) on the dataset.

Evaluation of various classification algorithms will be done by comparing evaluation metrics like the Accuracy, Precision and Area under the Curve (AUC) of each algorithm.

A conclusion is arrived at based on results.

Limitations / Challenges / Continuity

- Most of the Car feature had missing values and hence the feature could not be used. This is in-vehicle coupon recommendation, hence maybe Car feature was critical in determining whether the coupon would be accepted or not. After all, there are driving scenarios involved.
- This dataset is partially balanced. The results of supervised learning algorithms used to make predictions would skew slightly towards the class with the class with higher percentage of records. The percentage of accepted coupons: 56.843%. The percentage of rejected coupons: 43.157%. In the future, Synthetic Minority Over-sampling Technique (SMOTE) can be used to generate a more balanced dataset.
- The dataset was more focused on a particular type of population. The dataset should have been created by sampling all types of population. For example, Depari et al. (2022) found that the data contained mostly married females who like to travel alone on a sunny day around 6 PM. Most of them have attended college, yet didn't graduate (Depari et al., 2022). For those who have an occupation, it states that most of them earn an income of around \$25000 - \$37499 (Depari et al., 2022). It was also mentioned that the destination is mostly the No Urgent Place such as Coffee House, which provides a coupon that expires in one day (Depari et al., 2022).
- Patil et al. (2019) also observed that customers tend to purchase the same coupon over and over again. The dataset for in-vehicle coupon response is deficient in data over a periodic basis to help uncover such patterns. This is a limitation.

GitHub Repository

<https://github.com/suchetasikdar1/CIND820>

References

Ahmed, N., & Umair, M. (2024). *Churn prediction using machine learning: A coupon optimization technique*. In World Journal of Advanced Engineering Technology and Sciences (Vols. 12–12, Issue 02, pp. 332–354) [Journal-article].

<https://wjaets.com/sites/default/files/WJAETS-2024-0310.pdf>

Atiq, R., Fariha, F., Mahmud, M., Yeamin, S. S., Rushee, K. I., & Rahim, S. (2022, October 8). *A comparison of missing value imputation techniques on coupon acceptance prediction*. MECS PPress. <https://www.mecs-press.org/ijitcs/ijitcs-v14-n5/IJITCS-V14-N5-2.pdf>

Babaoglu C. (n.d.). CMTH 642 – Data Analytics: Advanced Methods course [PowerPoint slide printouts]

Babaoglu C. (2018, January 6). *How to conduct data analysis process Systematically* [Video]. YouTube. <https://www.youtube.com/watch?v=NMffLbFql5k>

Inyama C. (2023). *In-Vehicle-Coupon-Recommendations/coupon_Project (2).ipynb at main · chrisinyama/In-Vehicle-Coupon-Recommendations*. GitHub. [https://github.com/chrisinyama/In-Vehicle-Coupon-Recommendations/blob/main/coupon_Project%20\(2\).ipynb](https://github.com/chrisinyama/In-Vehicle-Coupon-Recommendations/blob/main/coupon_Project%20(2).ipynb)

Depari, G. S., Shu, E., Fachriza, C. A., Chow, J., Wijaya, J., & Winata, R. (2022). Customer's responses towards in-vehicle coupon recommendation an implementation of business analytics concept. *Jurnal Ekonomi* (Vol. 11, Issue 02) [Journal-article]

<https://download.garuda.kemdikbud.go.id/article.php?article=2994923&val=22616&title=CUSTOMER%20RESPONSES%20TOWARDS%20IN-VEHICLE%20COUPON%20RECOMMENDATION%20AN%20IMPLEMENTATION%20OF%20BUSINESS%20ANALYTICS%20CONCEPT>

Niralidedaniya. (2023, January 19). In Vehicle Coupon Recommendation — a Machine learning classification case study. *Medium*. <https://medium.com/@niralidedaniya/in-vehicle-coupon-recommendation-a-machine-learning-classification-case-study-df67e7835703>

Patil, Y., Pawar, O., & Ingle, D. R. (2019). Coupon Purchase Prediction using Machine Learning. In *IJSRD - International Journal for Scientific Research & Development* (Vol. 7, Issue 02) [Journal-article]. <https://ijsrd.com/articles/IJSRDV7I21146.pdf>

Supervised learning. (n.d.). Scikit-learn. https://scikit-learn.org/stable/supervised_learning.html

UCI Machine Learning Repository. (2020). <https://archive.ics.uci.edu/dataset/603/in+vehicle+coupon+recommendation>

Wang, T., Rudin, C., Doshi-Velez, F., Liu, Y., Klampfl, E., & MacNeille, P. (2017). A Bayesian framework for learning rule sets for interpretable classification. In Maya Gupta (Ed.), *Journal of Machine Learning Research* (Vols. 1–37) [Journal-article]. <https://jmlr.csail.mit.edu/papers/volume18/16-003/16-003.pdf>

How to run Machine Learning algorithms in GPU. (n.d.). Stack Overflow. <https://stackoverflow.com/questions/72985935/how-to-run-machine-learning-algorithms-in-gpu>

Ofir Shalev (@ofirdi). (2021, December 10). Recall, Precision, F1, ROC, AUC, and everything - The Startup - Medium. *Medium*. <https://medium.com/swlh/recall-precision-f1-roc-auc-and-everything-542aedef322b9>