

# Coupon Recommendation System Using User Behavior Data

HubbleMind



## Contents

.....	1
Introduction .....	3
Expected Outcome.....	3
Week 1: Data Understanding and Cleaning.....	3
Data cleaning and preprocessing.....	4
Week 2: Exploratory Data Analysis (EDA) .....	5
Week 3: Machine Learning Models .....	8
Train-Test Split: .....	8
Models performance: .....	9
<b>Random Forest Classifier</b> .....	11
Define the Parameter Grid:.....	13
Analysis review: .....	14

# Coupon Recommendation System Using User Behavior Data

---

## Introduction

As years have gone by in the era of commerce and retail market, customers base being one of the major reason in enhancing business makes competition increase among businesses. One of the most effective ways businesses stay on top of their game is by personalized discounts such as coupons and with how diverse humans have being over these years and their individual preference, providing them with the right coupon can be quite challenging and that's where machine learning (ML) can be key to understanding consumers behavior. This report shows each stage of the project from data preparation to model selection and hyperparameter tuning etc.

## Expected Outcome

- Predict which coupons a user is likely to find valuable based on their behavior.
- Increase coupon redemption rates and sales by personalizing offers.
- Provide businesses with insights into user preferences and behaviors, thereby improving marketing strategies.

This project aims to demonstrate the power of machine learning in solving practical, real-world problems in the retail and e-commerce domain, and to create a system that drives better customer experiences and business outcomes.

## Week 1: Data Understanding and Cleaning

The dataset utilized for this analysis, designated as in-vehicle-coupon-recommendation.csv encompasses multiple features to predict whether a user will accept a coupon based on several factors, such as weather, passenger information, time of day etc. The data comprises both numerical and categorical feature is detailed as follows:

### Numerical features:

- Temperature: The outside temperature in Fahrenheit.
- has\_children : ether the user has children or not.
- direction\_same: Whether the user is heading in the same direction as the coupon destination.
- Y (likely the target variable): (1 for accepting the coupon, 0 for rejecting it).
- toCoupon\_GEQ5min
- toCoupon\_GEQ15min
- toCoupon\_GEQ25min
- direction\_opp

### Categorical features:

- Destination: Where the user is going (e.g., No Urgent Place, Work).
- Passenger: Who the user is traveling with (e.g., Alone, Friend(s))

- Weather: The weather condition when the coupon was offered (e.g., Sunny, Rainy).
- Time: The time of day when the coupon was presented (e.g., 10AM, 2PM).
- Coupon: offered (e.g., Coffee House, Restaurant)
- car (has many null values)
- Bar (has some null values)
- Coffeehouse: (has some null values)
- CarryAway :(has some null values)
- RestaurantLessThan20 (has some null values)
- Restaurant20To50 (has some null values)
- Expiration:1day or 2 hours
- gender: male or female
- age
- marital Status
- education
- occupation
- income

## Data cleaning and preprocessing

This involves preparing raw data for analysis or modeling by improving its quality and structure. First, missing values are addressed by either removing columns/rows with excessive gaps or filling them with statistical measures (mean, median and mode). Duplicates are removed to ensure data uniqueness. Categorical data inconsistencies are resolved by standardizing entries. Features are encoded into numerical formats using methods like one-hot encoding. Numerical data is scaled using normalization or standardization for consistency.

**Import dataset:** Read the data from a file (e.g., CSV) into a Data Frame for processing and Identify which columns are **categorical** (e.g., weather, time) and which are **numerical** (e.g., temperature).

**Handle missing values:** this address missing data to ensure the dataset is clean and usable for analysis or modeling and it include:

**Drop columns with too many missing values:** Columns with over a certain percentage (e.g., 50%) of missing values are dropped as they may lack sufficient information.

**Handle Inconsistencies in Categorical Data:** reduces inconsistent labels can cause problems during analysis or modeling, especially with algorithms that rely on distinct values.

**Encode Categorical Variables:** Convert categorical variables into a format suitable for machine learning models by using **OneHotEncoder** to create binary (0/1) columns for each category in a feature.

**Scale numerical variables:** Many machine learning algorithms are sensitive to differences in scale, and standardized data improves performance thereby Normalizing numerical values to a standard range for improved model performance.

## Week 2: Exploratory Data Analysis (EDA)

This process of analyzing and visualizing data to uncover patterns, relationships, and insights. It involves summarizing data using descriptive statistics (mean, median and variance) and visualizations like histograms, scatter plots, and box plots. EDA helps identify trends, anomalies, missing values, and correlations, providing a deeper understanding of the dataset.

Visualize relationships between the target variable (Y) and categorical features like weather, time, passenger and age.

- In weather vs y, it shows that sunny weather has more positive results especially for Y=1 variable than Y=0.
- In time vs Y, it shows that 6pm to 7am has more positive outcome in both variables than other time value
- In passenger vs Y, people travelling with alone are more likely to accept the coupon or reject it than other group of people.
- For age, individuals from 21-30 has the highest outcome especially for Y=1

### Week 2:

#### Relationship within the data through visualization and correlation analysis

```
In [113.. #Relationship between target variable(y) and categorical features such as weather, time, passenger and age
```

```
In [114.. import seaborn as sns
import matplotlib.pyplot as plt
```

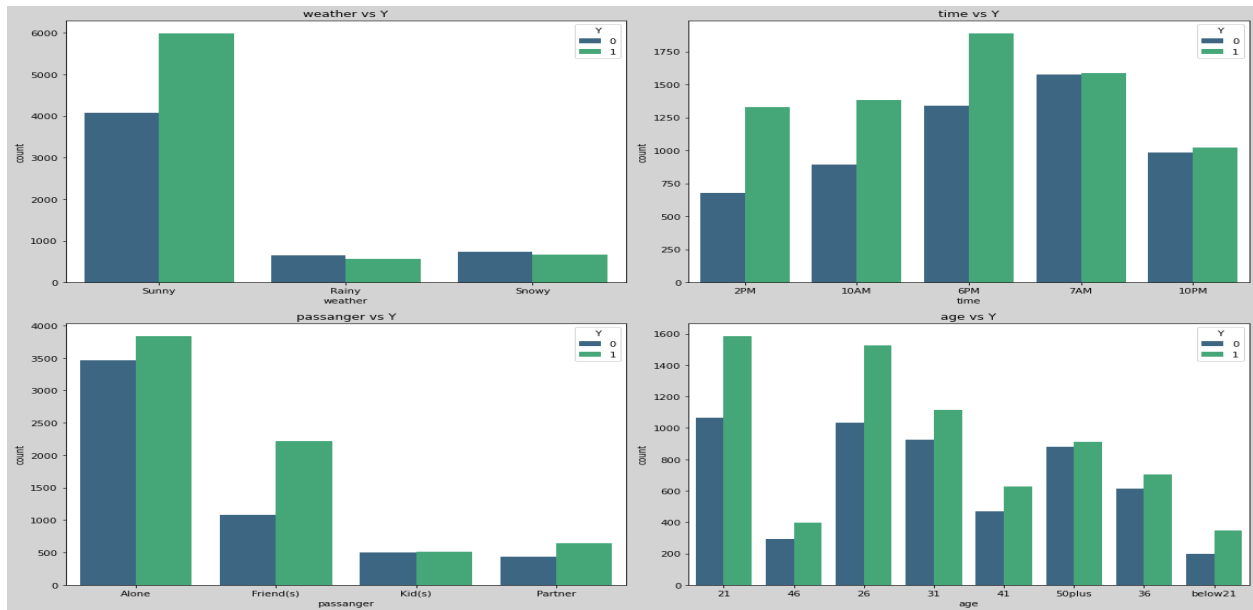
```
In [115.. def plot_categorical_vs_target(vehicle, categorical_features, target):
fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(16, 12))
axes = axes.flatten()
for i, feature in enumerate(categorical_features):
sns.countplot(data=vehicle, x=feature, hue=target, ax=axes[i], palette="viridis")
axes[i].set_title(f'{feature} vs {target}')
axes[i].legend(title=target, loc='upper right')
plt.tight_layout()
plt.show()

categorical_features = ['weather', 'time', 'passanger', 'age']
plot_categorical_vs_target(vehicle, categorical_features, 'Y')
```

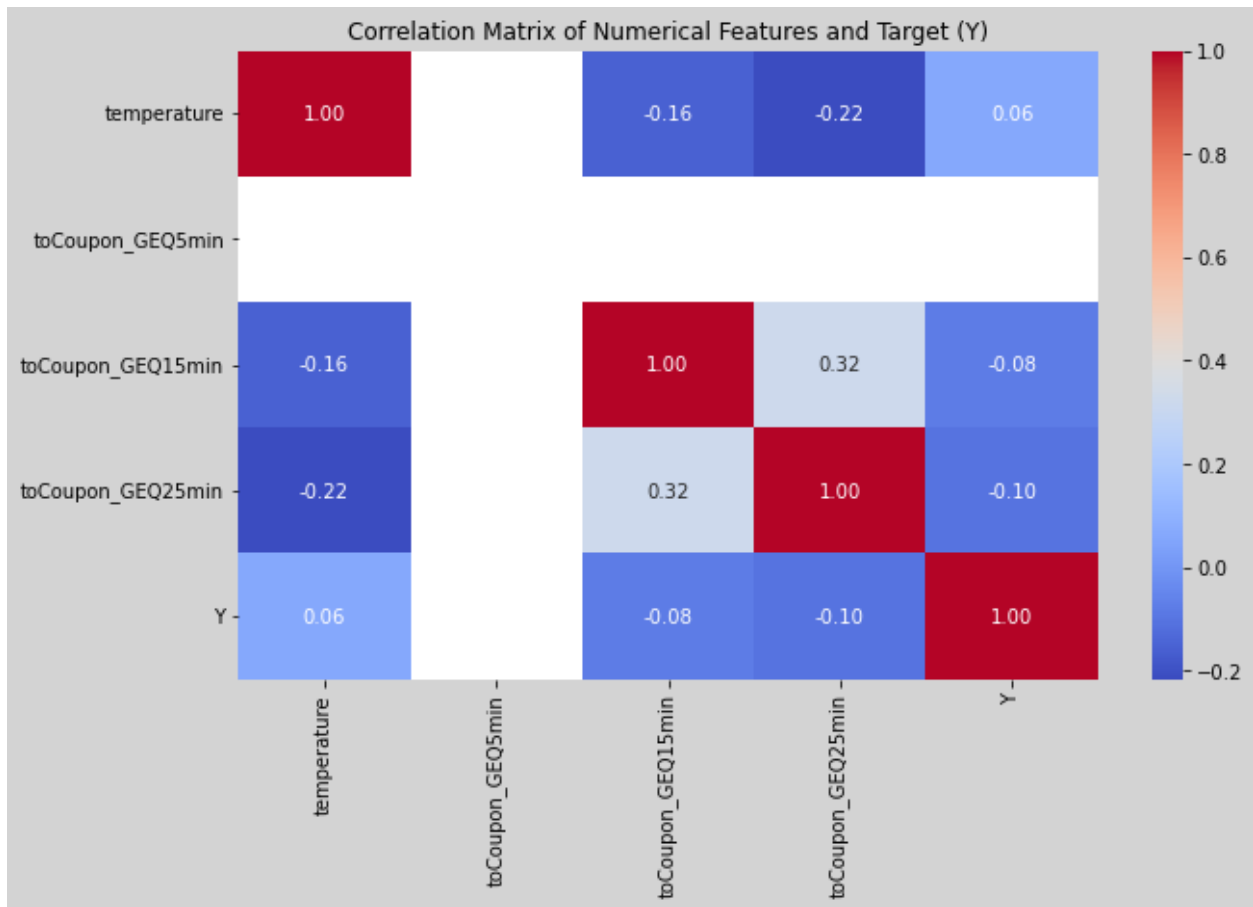
weather vs Y

time vs Y

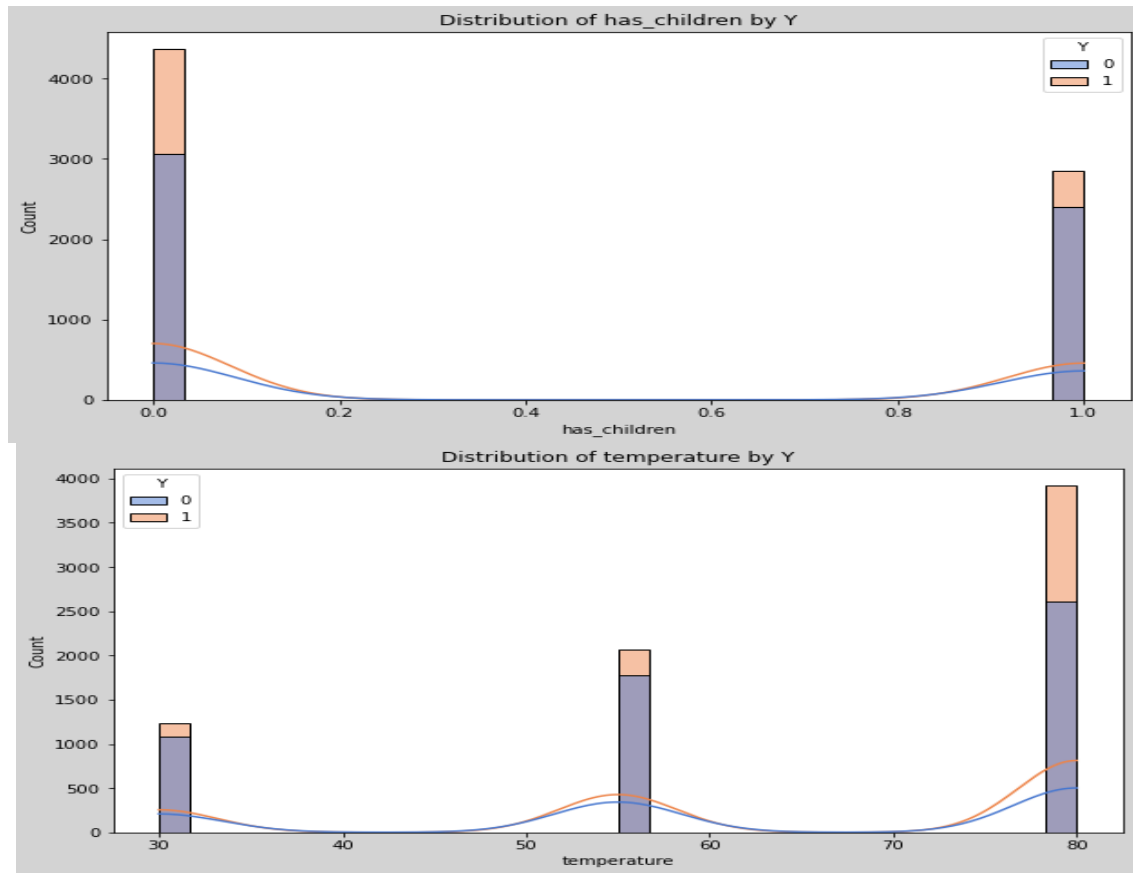
Activate Windows  
Go to Settings to activate Windows



- Correlation analysis between numerical features and the target variable (Y)



- Visualize distributions of key numerical features and their impact on coupon acceptance.



1 for accepting the coupon, 0 for rejecting it therefore this shows that more people rejected the coupon especially when the temperature is at an all-time high.

For this particular graph, in Y vs children, they are more likely to receive or reject the coupon when there are no children available

## Week 3: Machine Learning Models

A **machine learning model** is a mathematical representation or algorithm designed to learn patterns from data. These models enable computers to make predictions, decisions, or categorizations without being explicitly programmed for every specific task. Common ML models such as:

- Linear Regression: Predicts a continuous output using a linear relationship.
- Logistic Regression: Classifies binary or multi-class data.
- Decision Trees: Splits data based on conditions to make decisions.
- Random Forests: Combines multiple decision trees for improved accuracy.
- Support Vector Machines (SVM): Finds a boundary to separate classes in a dataset.
- Neural Networks: Mimics the human brain to model complex relationships.

### Train-Test Split:

The **train-test split** is a method to evaluate machine learning models by dividing the dataset into two parts: the training set (70–80%) and the test set (20–30%). The model learns patterns from the training set and is then evaluated on the test set to check its performance on unseen data. This approach ensures the model can generalize well and prevents over fitting, where the model memorizes the training data but performs poorly on new data. For imbalanced datasets, stratified splitting is recommended to maintain class proportions. It's a simple yet essential step for building robust and reliable machine learning models. Example of the model is:

[illegible]



```

In [122.. from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=100)

In [123.. #80 % percent training
x_train

Out[123.. destination passenger weather temperature time coupon expiration gender age maritalStatus ... Work Rainy Snowy Sunny Home No Urgent Place Work Rainy Snowy Sunny
710 1 3 2 80 3 4 1 1 3 3 ... 0.0 0.0 0.0 1.0 0.0 1.0 0.0 0.0 0.0 1.0
12015 0 3 0 55 3 1 0 1 0 2 ... 0.0 1.0 0.0 0.0 1.0 0.0 0.0 1.0 0.0 0.0
5864 0 0 2 80 3 4 1 0 1 1 ... 0.0 0.0 0.0 1.0 1.0 0.0 0.0 0.0 0.0 1.0
5656 1 1 2 80 3 4 0 0 2 2 ... 0.0 0.0 0.0 1.0 0.0 1.0 0.0 0.0 0.0 1.0
3472 1 1 2 80 0 1 1 0 4 1 ... 0.0 0.0 0.0 1.0 0.0 1.0 0.0 0.0 0.0 1.0
... ..
79 0 0 2 55 3 0 0 1 5 1 ... 0.0 0.0 0.0 1.0 1.0 0.0 0.0 0.0 0.0 1.0
12119 0 0 2 30 3 1 0 0 3 1 ... 0.0 0.0 0.0 1.0 1.0 0.0 0.0 0.0 0.0 1.0
8039 1 1 2 80 2 3 1 1 2 2 ... 0.0 0.0 0.0 1.0 0.0 1.0 0.0 0.0 0.0 1.0
6936 1 1 2 80 2 3 1 0 4 1 ... 0.0 0.0 0.0 1.0 0.0 1.0 0.0 0.0 0.0 1.0
5640 0 0 2 55 1 1 0 0 3 2 ... 0.0 0.0 0.0 1.0 1.0 0.0 0.0 0.0 0.0 1.0

10147 rows x 37 columns

In [124.. # 20% testing set.
x_test

```

# Models performance:

Looking at code from my jupyter notebook on Logistic Regression, Decision Tree Classifier and Random Forest classifier.

## LogisticRegression

```

In [125.. from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report

In [126.. model=LogisticRegression()

In [127.. model.fit(x_train,y_train)

C:\Users\DELL\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:763: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_i = _check_optimize_result(
LogisticRegression()

Out[127..

In [128.. y_pred=model.predict(x_test)

```

```

In [129.. model.score(x_test,y_test)

```

```

Out[129.. 0.6342136381553015

```

```

In [130.. from sklearn.metrics import confusion_matrix

```

```

In [131.. print(confusion_matrix(y_test,y_pred))

```

```

[[ 517 601]
 [ 327 1092]]

```

```

In [132.. from sklearn.metrics import classification_report

```

```

In [133.. print(classification_report(y_test,y_pred))

```

	precision	recall	f1-score	support
0	0.61	0.46	0.53	1118
1	0.65	0.77	0.70	1419
accuracy			0.63	2537
macro avg	0.63	0.62	0.61	2537
weighted avg	0.63	0.63	0.62	2537

# DecisionTreeClassifier

```
In [134... from sklearn.tree import DecisionTreeClassifier, export_text, plot_tree
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```

```
In [135... dtc=DecisionTreeClassifier()
```

```
In [136... dtc.get_params()
```

```
Out[136... {'ccp_alpha': 0.0,
'class_weight': None,
'criterion': 'gini',
'max_depth': None,
'max_features': None,
'max_leaf_nodes': None,
'min_impurity_decrease': 0.0,
'min_impurity_split': None,
'min_samples_leaf': 1,
'min_samples_split': 2,
'min_weight_fraction_leaf': 0.0,
'random_state': None,
'splitter': 'best'}
```

```
In [137... dtc.fit(x_train,y_train)
```

```
min_samples_split': 2,
'min_weight_fraction_leaf': 0.0,
'random_state': None,
'splitter': 'best'}
```

```
In [137... dtc.fit(x_train,y_train)
```

```
Out[137... DecisionTreeClassifier()
```

```
In [138... Y_pred=dtc.predict(x_test)
```

```
In [139... print(confusion_matrix(y_test,Y_pred))
```

```
[[ 701  417]
 [ 417 1002]]
```

```
In [140... print(classification_report(y_test,Y_pred))
```

	precision	recall	f1-score	support
0	0.63	0.63	0.63	1118
1	0.71	0.71	0.71	1419
accuracy			0.67	2537
macro avg	0.67	0.67	0.67	2537
weighted avg	0.67	0.67	0.67	2537

## Random Forest Classifier

```
In [153... from sklearn.ensemble import RandomForestClassifier
```

```
In [154... rf = RandomForestClassifier()
```

```
In [155... rf.fit(x_train,y_train)
```

```
Out[155... RandomForestClassifier()
```

```
In [156... y1_pred=rf.predict(x_test)
```

```
In [157... #accuracy  
rf.score(x_test,y_test)
```

```
Out[157... 0.7299960583366181
```

```
In [158... accuracy = accuracy_score(y_test, y1_pred)  
accuracy
```

```
Out[158... 0.7299960583366181
```

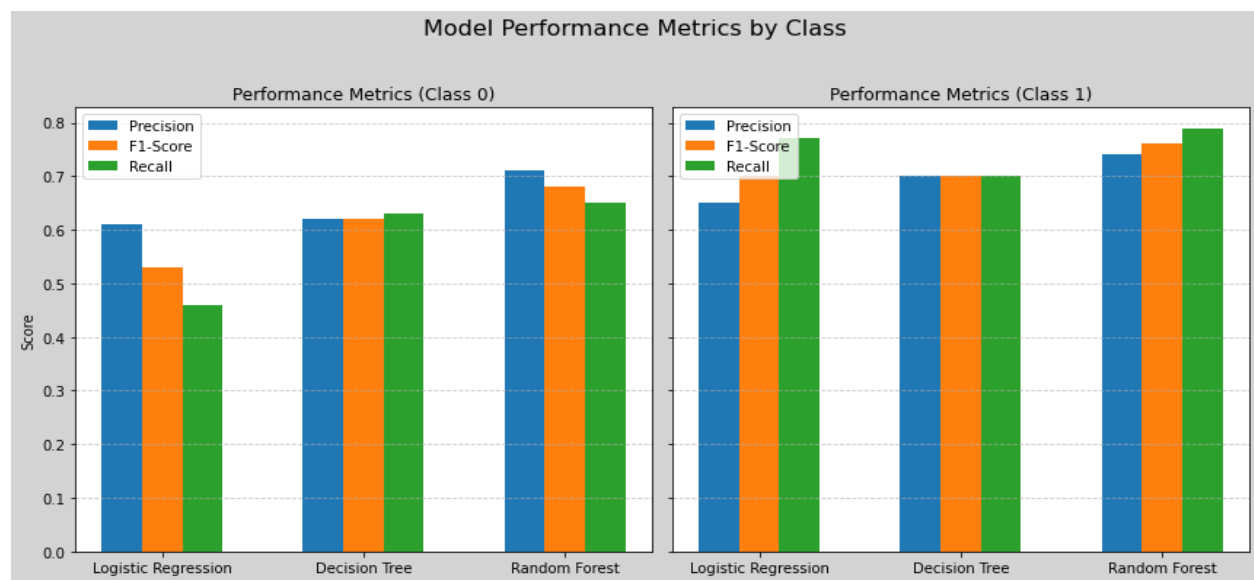
```
In [159... from sklearn.metrics import accuracy_score, confusion_matrix, classification_report  
import random
```

```
In [160... print(classification_report(y_test,y1_pred))
```

	precision	recall	f1-score	support
0	0.71	0.66	0.68	1118
1	0.74	0.79	0.77	1419
accuracy			0.73	2537
macro avg	0.73	0.72	0.72	2537
weighted avg	0.73	0.73	0.73	2537

1 for accepting the coupon, 0 for rejecting it

Models	Y	accuracy	precision	F1-Score	Recall	support
Logistic Regression	0	0.63	0.61	0.53	0.46	1118
	1		0.65	0.70	0.77	1419
Decision tree classifier	0	0.67	0.62	0.62	0.63	1118
	1		0.70	0.70	0.70	1419
Random Forest Classifier	0	0.73	0.71	0.68	0.65	1118
	1		0.74	0.76	0.79	1419



It clearly shows that Random forest model for Y variable 0 and 1 has the superior and best performance model in comparison with the different models tested. It consistently has the highest precision, recall, and F1-Score, making it the most reliable and effective choice among the three models. Random Forest's ensemble approach (using multiple decision trees) enhances its ability to generalize, reducing over fitting and ensuring stable performance on unseen data and its ability to maintain strong performance across all metrics for both classes hence making it suitable.

## Week 4: Fine-Tuning and Reporting

Random Forest, an ensemble learning method, has several hyper parameters that influence its performance, such as the number of trees (`n_estimators`), maximum tree depth (`max_depth`), and the minimum samples required to split a node (`min_samples_split`). This involves optimizing its hyper parameters to achieve the best performance. Cross-validation splits the training dataset into several folds, using each fold as validation while training on the rest, ensuring reliable performance estimates.

### Define the Parameter Grid:

Create a dictionary of hyperparameters to search. Common hyper parameters for Random Forest include:

- `n_estimators`: Number of trees in the forest.
- `max_depth`: Maximum depth of each tree.
- `min_samples_split`: Minimum number of samples required to split an internal node.
- `min_samples_leaf`: Minimum number of samples required at a leaf node.
- `max_features`: Number of features to consider when looking for the best split.

```
In [163.. from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV, train_test_split

In [164.. #Number of trees in random forest
n_estimators=[int(x) for x in np.linspace(start=10,stop=80,num=10)]
#Number of features to consider at every split
max_features=['auto','sqrt']
#maximum number of levels in tree
max_depth=[2,4]
#maximum number of samples required to split a node
min_samples_split=[2,5]
#minimum number of samples required at each leaf node
min_samples_leaf=[1,2]
#Method of selecting samples of training each tree
bootstrap=[True,False]

In [165.. param_grid={'n_estimators': n_estimators,
                    'max_features': max_features,
                    'max_depth': max_depth,
                    'min_samples_split': min_samples_split,
                    'min_samples_leaf': min_samples_leaf,
                    'bootstrap': bootstrap}

print(param_grid)

{'n_estimators': [10, 17, 25, 33, 41, 48, 56, 64, 72, 80], 'max_features': ['auto', 'sqrt'], 'max_depth': [2, 4], 'min_samples_split': [2, 5], 'min_samples_
leaf': [1, 2], 'bootstrap': [True, False]}

In [166.. rf=RandomForestClassifier()
```

```

In [167]: rf_Grid=GridSearchCV(estimator=rf,param_grid=param_grid ,cv=3,verbose=2,n_jobs=4)

In [168]: rf_Grid.fit(x_train,y_train)

Out[168]: Fitting 3 folds for each of 320 candidates, totalling 960 fits
GridSearchCV(cv=3, estimator=RandomForestClassifier(), n_jobs=4,
             param_grid={'bootstrap': [True, False], 'max_depth': [2, 4],
                          'max_features': ['auto', 'sqrt'],
                          'min_samples_leaf': [1, 2],
                          'min_samples_split': [2, 5],
                          'n_estimators': [10, 17, 25, 33, 41, 48, 56, 64, 72,
                                             80]},
             verbose=2)

In [169]: rf_Grid.best_params_

Out[169]: {'bootstrap': False,
           'max_depth': 4,
           'max_features': 'sqrt',
           'min_samples_leaf': 2,
           'min_samples_split': 5,
           'n_estimators': 56}

In [170]: print(f'Train Accuracy -: {rf_Grid.score(x_train,y_train):.3f}')
           print(f'Test Accuracy -: {rf_Grid.score(x_test,y_test):.3f}')

Train Accuracy -: 0.655
Test Accuracy -: 0.644

```

Activate Windows  
Go to Settings to activate Windows.

## Analysis review:

From the analysis done, the Random Forest model was fine-tuned with GridSearchCV. These settings constrain tree growth (max\_depth=4), limit overfitting by requiring a minimum of 2 samples per leaf, and use the square root of features for splits. The model uses 56 trees without bootstrapping (bootstrap=False), encouraging diversity. This configuration balances bias and variance, improving generalization. After training, the optimized model is evaluated using cross-validation and achieves robust performance. Final metrics confirm improved predictive accuracy and reliability on test data.

The results after fine-tuning the Random Forest model using GridSearchCV show a **Train Accuracy of 0.655** and a **Test Accuracy of 0.64**. These values indicate that the model is generalizing reasonably well, with a small gap between training and testing accuracy. Consequently, this model may be confidently employed to predict productivity with enhanced accuracy and reliability, establishing it as an effective solution for our application.

## Python Code from my project

<https://github.com/Chizobaeze/Machine-Learning-project-on-Coupon-Recommendation-System-Using-User-Behavior-Data/blob/main/Hubblemind%20Eze%20chizoba.ipynb>