



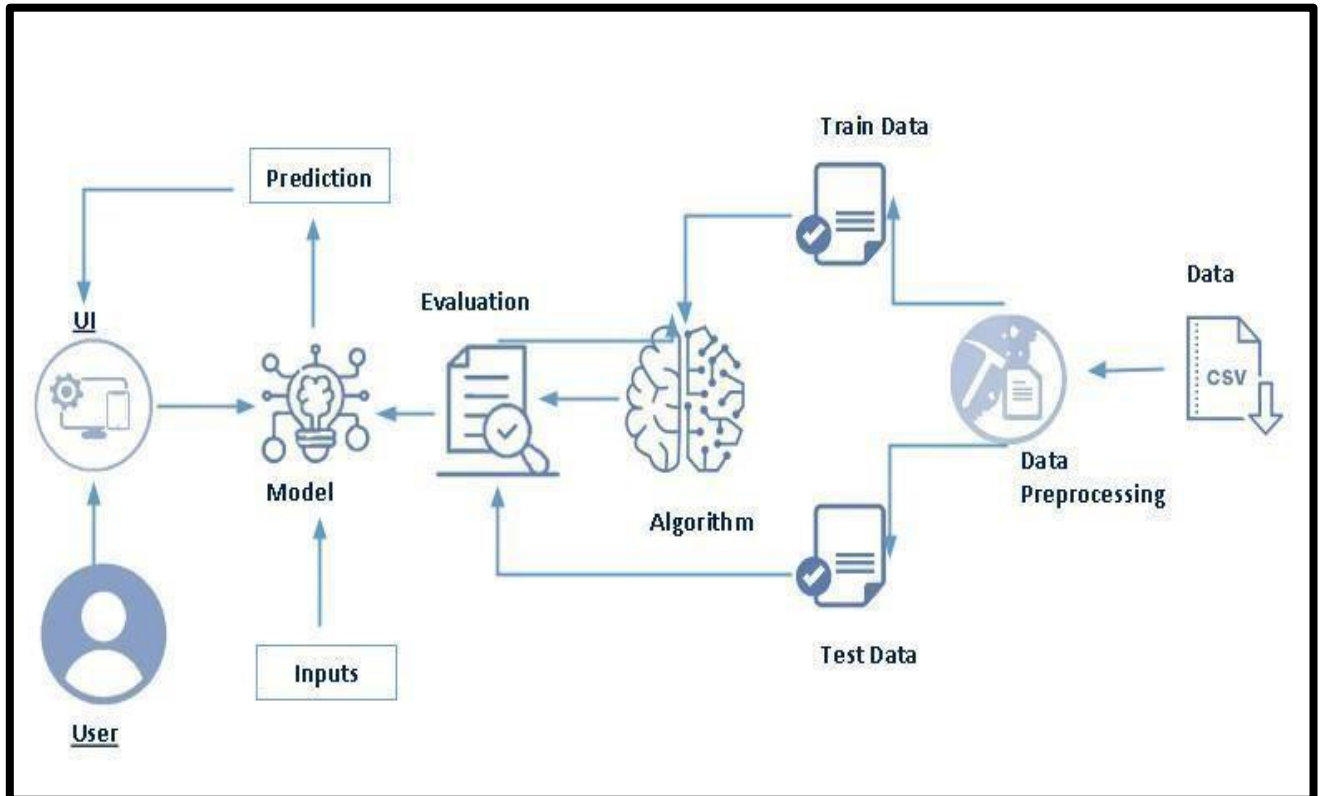
ECOMMERCE SHIPPING PREDICTION USING MACHINE LEARNING

**BY- ARJUN KADTANE, YAGYESH KUMAR, NITIN VOONA, CHIDVILISA
YEPURI**

Ecommerce Shipping Prediction Using Machine Learning

Ecommerce shipping prediction is the process of estimating the whether the product reached on time. which is based on various factors such as the origin and destination of the package, the shipping method selected by the customer, the carrier used for shipping, and any potential delays or issues that may arise during the shipping process. Machine learning models can be used to make accurate predictions about shipping times based on historical data and real-time updates from carriers. These models may take into account factors such as weather conditions, traffic, and other external factors that can impact delivery times. Over All Ecommerce shipping prediction is an important tool for ecommerce businesses that want to provide accurate delivery estimates to their customers and improve their overall customer experience.

Technical Architecture:



Project Flow:

- User interacts with the UI to enter the input.
- Entered input is analysed by the model which is integrated.
- Once model analyses the input the prediction is showcased on the UI

To accomplish this, we have to complete all the activities listed below,

- Define Problem / Problem Understanding
 - Specify the business problem
 - Business requirements
 - Literature Survey
 - Social or Business Impact.
- Data Collection & Preparation
 - Collect the dataset
 - Data Preparation
- Exploratory Data Analysis
 - Descriptive statistical
 - Visual Analysis
- Model Building
 - Training the model in multiple algorithms
 - Testing the model
- Performance Testing & Hyperparameter Tuning
 - Testing model with multiple evaluation metrics
 - Comparing model accuracy before & after applying hyperparameter tuning
- Model Deployment
 - Save the best model
 - Integrate with Web Framework
- Project Demonstration & Documentation
 - Record explanation Video for project end to end solution
 - Project Documentation-Step by step project development procedure

Prior Knowledge:

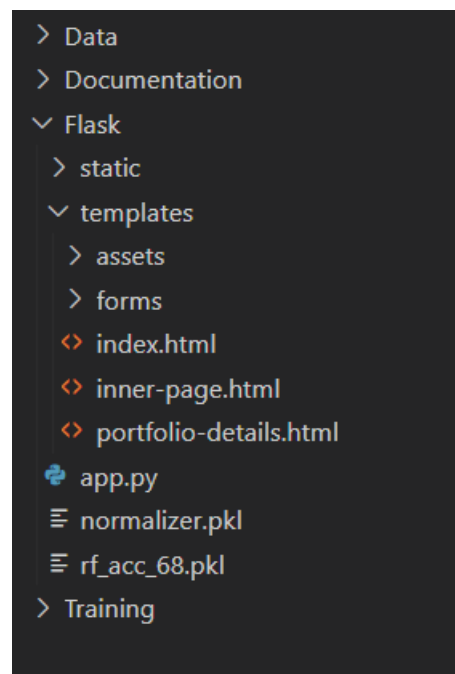
You must have prior knowledge of following topics to complete this project.

ML Concepts

- Supervised learning: <https://www.javatpoint.com/supervised-machine-learning>
- Unsupervised learning: <https://www.javatpoint.com/unsupervised-machine-learning>
- Decision tree: <https://www.javatpoint.com/machine-learning-decision-tree-classificationalgorithm>
- Random forest: <https://www.javatpoint.com/machine-learning-random-forest-algorithm>
- KNN: <https://www.javatpoint.com/k-nearest-neighbor-algorithm-for-machine-learning>
- Xgboost: <https://www.analyticsvidhya.com/blog/2018/09/an-end-to-end-guide-to-understand-the-math-behind-xgboost/>
- Evaluation metrics: <https://www.analyticsvidhya.com/blog/2019/08/11-important-model-evaluation-error-metrics/>
- Flask Basics : https://www.youtube.com/watch?v=lj4I_CvBnt0

Project Structure:

Create the Project folder which contains files as shown below



- We are building a flask application which needs HTML pages stored in the templates folder and a python script app.py for scripting.
- Static folder and assets folder contain the CSS and JavaScript files along with images.
- rf_acc_68.pkl and normalizer.pkl are our saved models. Further we will use these models for flask integration.
- Training folder contains a model training file that are jupyter notebooks.

Milestone 1: Define Problem / Problem Understanding

Activity 1: Specify the business problem

Refer Project Description

Activity 2: Business requirements

Here are some potential business requirements for an ecommerce product delivery estimation predictor using machine learning:

1. **Accurate Delivery Estimates:** The primary goal of the system should be to provide accurate delivery estimates to customers. The estimated delivery time should consider factors such as distance, traffic, weather, and any other relevant variables.
2. **Real-Time Updates:** Customers should be able to receive real-time updates on the status of their delivery, including any delays or changes to the estimated delivery time. The system should be able to adjust the estimated delivery time based on the most up-to-date information available.
3. **Integration with Ecommerce Platforms:** The system should be able to integrate with popular ecommerce platforms, such as Shopify or Magento, to automatically retrieve order information and provide accurate delivery estimates.
4. **Machine Learning Models:** The system should use machine learning models to predict delivery times based on historical delivery data and other relevant variables. The machine learning models should be trained and optimized to improve accuracy over time.
5. **Scalability:** The system should be scalable to handle large volumes of orders and delivery estimates. It should be able to calculate delivery estimates quickly and accurately for a high number of orders at the same time.
6. **Reporting and Analytics:** The system should provide reporting and analytics capabilities to help ecommerce businesses track delivery performance and identify areas for improvement.
7. **Customer Service Integration:** The system should be integrated with customer service channels, such as email or chat support, to provide timely updates to customers and handle any delivery-related inquiries or issues.
8. **Cost-Effective:** The system should be cost-effective to implement and maintain, with reasonable pricing models and minimal upfront costs.

Activity 3: Literature Survey (Student Will Write)

A literature survey for a Ecommerce Shipping Prediction project would involve researching and reviewing existing studies, articles, and other publications on the topic of Ecommerce Shipping Prediction. The survey would aim to gather information on current systems, their strengths and weaknesses, and any gaps in knowledge that the project could address. The literature survey would also look at the methods and techniques used in previous projects, and any relevant data or findings that could inform the design and implementation of the current project.

Activity 4: Social or Business Impact.

Social Impacts:

1. **Improved Customer Experience:** Providing accurate delivery estimates can help improve the overall customer experience by reducing uncertainty and increasing transparency.
2. **Reduced Environmental Impact:** By accurately estimating delivery times, businesses can optimize their logistics and reduce the number of unnecessary trips and emissions from delivery vehicles.
3. **Reduced Stress for Delivery Workers:** By optimizing delivery routes and times, businesses can reduce the workload and stress on delivery workers, leading to a better work environment and potentially reducing turnover.

Business Impacts:

1. **Increased Sales:** Accurate delivery estimates can help increase customer confidence and reduce cart abandonment rates, leading to increased sales and revenue.
2. **Improved Operational Efficiency:** By optimizing delivery routes and times, businesses can reduce costs associated with transportation, labor, and inventory management.
3. **Competitive Advantage:** Implementing a delivery estimation predictor using machine learning can provide a competitive advantage over other ecommerce businesses that do not offer accurate and transparent delivery estimates.
4. **Better Data-Driven Decision Making:** By analyzing delivery data and performance metrics, businesses can make data-driven decisions to optimize their logistics and improve their overall delivery performance.
5. **Brand Loyalty:** Providing accurate delivery estimates and real-time updates can help build brand loyalty by increasing trust and confidence in the business.

Milestone 2: Data Collection & Preparation

ML depends heavily on data. It is the most crucial aspect that makes algorithm training possible. So, this section allows you to download the required dataset.

Activity 1: Collect the dataset.

There are many popular open sources for collecting the data. Eg: kaggle.com, UCI repository, etc.

In this project we have used .csv data. This data is downloaded from kaggle.com. Please refer to the link given below to download the dataset.

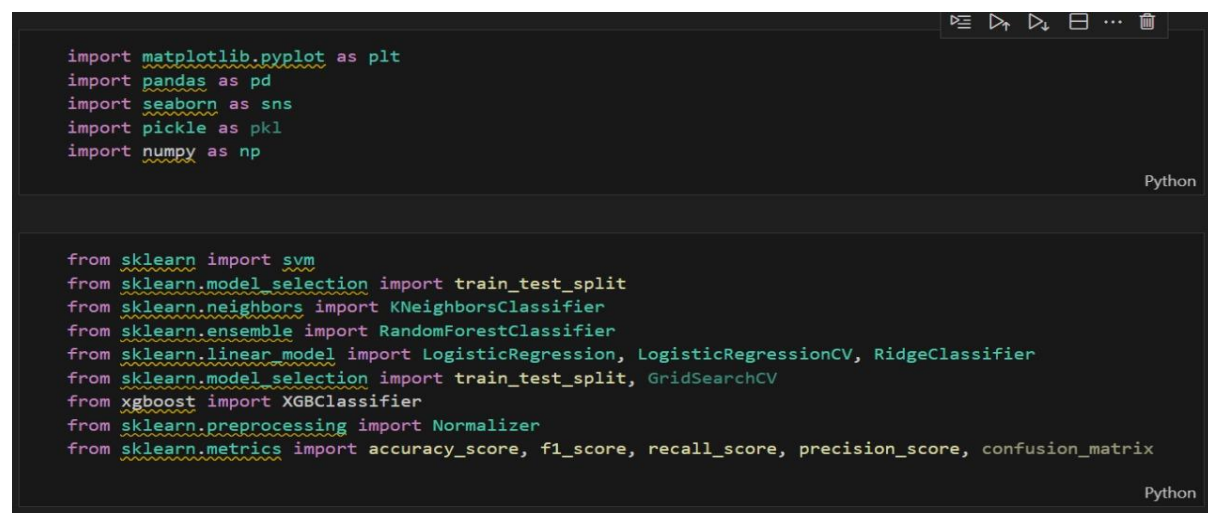
Link: <https://www.kaggle.com/datasets/prachi13/customer-analytics?select=Train.csv>

As the dataset is downloaded. Let us read and understand the data properly with the help of some visualisation techniques and some analysing techniques.

Note: There are a number of techniques for understanding the data. But here we have used some of it. In an additional way, you can use multiple techniques.

Activity 1.1: Importing the libraries

Import the necessary libraries as shown in the image.



```
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
import pickle as pkl
import numpy as np

from sklearn import svm
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression, LogisticRegressionCV, RidgeClassifier
from sklearn.model_selection import train_test_split, GridSearchCV
from xgboost import XGBClassifier
from sklearn.preprocessing import Normalizer
from sklearn.metrics import accuracy_score, f1_score, recall_score, precision_score, confusion_matrix
```

Activity 1.2: Read the Dataset

Our dataset format might be in .csv, excel files, .txt, .json, etc. We can read the dataset with the help of pandas.

In pandas we have a function called read_csv() to read the dataset. As a parameter we have to give the directory of the csv file.

```
data = pd.read_csv("Train.csv")
```

Python

```
data.head()
```

Python

	ID	Warehouse_block	Mode_of_Shipment	Customer_care_calls	Customer_rating	Cost_of_the_Product	Prior_purcha
0	1	D	Flight	4	2	177	
1	2	F	Flight	4	5	216	
2	3	A	Flight	2	2	183	
3	4	B	Flight	3	3	176	
4	5	C	Flight	2	2	184	

Activity 2: Data Preparation

As we have understood how the data is, let's pre-process the collected data.

The download data set is not suitable for training the machine learning model as it might have so much randomness so we need to clean the dataset properly in order to fetch good results.

This activity includes the following steps.

- Handling missing values
- Handling categorical data
- Handling Outliers

Note: These are the general steps of pre-processing the data before using it for machine learning. Depending on the condition of your dataset, you may or may not have to go through all these steps.

Activity 2.1: Handling missing values

- Let's find the shape of our dataset first. To find the shape of our data, the `df.shape` method is used. To find the data type, `df.info()` function is used.


```
data.shape

(10999, 12)

data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10999 entries, 0 to 10998
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   ID                     10999 non-null  int64
1   Warehouse_block        10999 non-null  object
2   Mode_of_Shipment        10999 non-null  object
3   Customer_care_calls     10999 non-null  int64
4   Customer_rating         10999 non-null  int64
5   Cost_of_the_Product     10999 non-null  int64
6   Prior_purchases         10999 non-null  int64
7   Product_importance      10999 non-null  object
8   Gender                  10999 non-null  object
9   Discount_offered        10999 non-null  int64
10  Weight_in_gms           10999 non-null  int64
11  Reached.on.Time_Y.N     10999 non-null  int64
dtypes: int64(8), object(4)
memory usage: 1.0+ MB
```

- For checking the null values, `df.isnull()` function is used. To sum those null values we use `.sum()` function. From the below image we found that there are no null values present in our dataset. So we can skip handling the missing values step.

```
data.isnull().sum()

ID                0
Warehouse_block    0
Mode_of_Shipment   0
Customer_care_calls 0
Customer_rating    0
Cost_of_the_Product 0
Prior_purchases    0
Product_importance 0
Gender             0
Discount_offered    0
Weight_in_gms       0
Reached.on.Time_Y.N 0
dtype: int64
```

Activity 2.2: Handling Categorical Values

As we can see our dataset has categorical data we must convert the categorical data to integer encoding or binary encoding.

To convert the categorical features into numerical features we use encoding techniques. There are several techniques but in our project we are using manual encoding with the help of list comprehension.

- In our project, categorical features are
 - Warehouse_block
 - Mode_of_shipment
 - Product_importance
 - Gender.

With list comprehension encoding is done.

```
label_map={}

for i in data.columns:
    if str(data[i].dtype) == 'object':
        temp={}
        cats=data[i].unique()
        for index in range(len(cats)):
            temp[cats [index]]=index

        label_map[i]=temp
        #Labeling
        data[i]=data[i].map(temp)
label_map

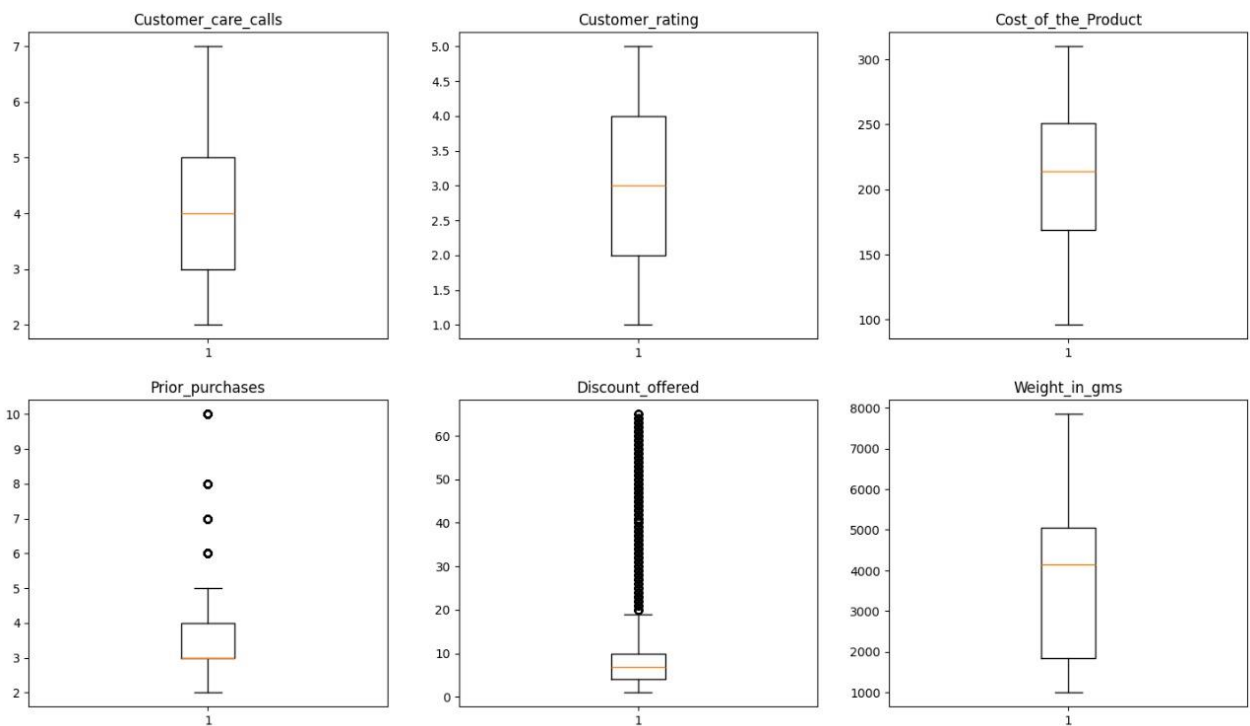
{'Warehouse_block': {'D': 0, 'F': 1, 'A': 2, 'B': 3, 'C': 4},
'Mode_of_Shipment': {'Flight': 0, 'Ship': 1, 'Road': 2},
'Product_importance': {'low': 0, 'medium': 1, 'high': 2},
'Gender': {'F': 0, 'M': 1}}
```

Activity 2.3: Handling Outliers in Data

With the help of boxplot, outliers are visualized. And here we are going to find upper bound and lower bound of numerical features with some mathematical formula.

- From the below diagram, we could visualize that Discount_offered, Prior_purchases features have outliers. Boxplot from matplotlib library is used here.

```
c=0
plt.figure(figsize=(18, 10))
for i in data.drop(columns=[ 'Warehouse_block', 'Mode_of_Shipment', 'Product_importance', 'Gender', 'Reach
    if str(data[i].dtype)=='object':
        continue
    plt.subplot(2, 3, c+1)
    plt.boxplot(data[i])
    plt.title(i)
    c+=1
plt.show()
```



- To find upper bound we have to multiply IQR (Interquartile range) with 1.5 and add it with 3rd quantile. To find lower bound instead of adding, subtract it with 1st quantile. Take image attached below as your reference.

```
def check_outliers(arr):
    Q1 = np.percentile(arr, 25, interpolation = 'midpoint')
    Q3 = np.percentile(arr, 75, interpolation = 'midpoint')
    IQR = Q3 - Q1
    #Above Upper bound
    upper = Q3 + 1.5 * IQR
    upper_array = np.array(arr)[arr > upper]
    print(' *3, len(upper_array[upper_array == True]), 'are over the upper bound:', upper)
    #Below Lower bound
    lower = Q1 - 1.5 * IQR
    lower_array = np.array(arr)[arr < lower]
    print(' *3, len(lower_array[lower_array == True]), 'are less than the lower bound:', lower, '\n')
    for i in data.drop(columns=[
        'Warehouse_block', 'Mode_of_Shipment', 'Product_importance', 'Gender', 'Reached.on.Time_Y.N', 'ID'
    ]).columns:
        if str(data[i].dtype) == 'object':
            continue
        print(i)
        check_outliers(data[i])

Customer_care_calls
0 are over the upper bound: 8.0
0 are less than the lower bound: 0.0

Customer_rating
0 are over the upper bound: 7.0
0 are less than the lower bound: -1.0

Cost_of_the_Product
0 are over the upper bound: 374.0
0 are less than the lower bound: 46.0

Prior_purchases
1003 are over the upper bound: 5.5
0 are less than the lower bound: 1.5

Discount_offered
2262 are over the upper bound: 19.0
0 are less than the lower bound: -5.0

Weight_in_gms
0 are over the upper bound: 9865.75
0 are less than the lower bound: -2976.25

C:\Users\Arjun\Anaconda3\envs\env\python.exe: DeprecationWarning: the 'interpolation=' argument to percentile was renamed to 'method=', which has additional options.
Users of the modes 'nearest', 'lower', 'higher', or 'midpoint' are encouraged to review the method they used. (Deprecated NumPy 1.22)
check_outliers(data[i])
```

- To handle the outliers transformation technique is used. Here L1 transformation is used.

➔ Data splitting

The data was split into train and test variables as shown below using the `train_test_split()` method of `scikitlearn` module with a `split_size` of 0.20 and a `random_state = 1234`.

```
x_train,x_test,y_train,y_test=train_test_split(data.drop(columns=['ID', 'Reached.on.Time_Y.N']),data['Reac
print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)
```

(8799, 10)
(2200, 10)
(8799,)
(2200,)

➔ Normalization

The data will be normalized using L1 regularisation that will be applied on `x_train` and `x_test` variables separately.

Milestone 3: Exploratory Data Analysis

Activity 1: Descriptive statistical

Descriptive analysis is to study the basic features of data with the statistical process. Here pandas has a worthy function called `describe`. With this `describe` function we can understand the unique, top and frequent values of categorical features. And we can find mean, std, min, max and percentile values of continuous features.

Codiumate: Options | Test this function
`data.describe(include='all')`

	ID	Warehouse_block	Mode_of_Shipment	Customer_care_calls	Customer_rating	Cost_of_the_Product
count	10999.00000	10999.000000	10999.000000	10999.000000	10999.000000	10999.000000
mean	5500.00000	1.833167	0.998454	4.054459	2.990545	210.196836
std	3175.28214	1.343823	0.567099	1.141490	1.413603	48.063272
min	1.00000	0.000000	0.000000	2.000000	1.000000	96.000000
25%	2750.50000	1.000000	1.000000	3.000000	2.000000	169.000000
50%	5500.00000	1.000000	1.000000	4.000000	3.000000	214.000000
75%	8249.50000	3.000000	1.000000	5.000000	4.000000	251.000000
max	10999.00000	4.000000	2.000000	7.000000	5.000000	310.000000

Activity 2: Visual analysis

Visual analysis is the process of using visual representations, such as charts, plots, and graphs, to explore and understand data. It is a way to quickly identify patterns, trends, and outliers in the data, which can help to gain insights and make informed decisions.

Activity 2.1: Univariate analysis

In simple words, univariate analysis is understanding the data with single feature. Here we have displayed two different graphs such as `histplot` and `countplot`.

Seaborn package provides a wonderful function `histplot`. With the help of `histplot`, we can find the distribution of the feature. To make multiple graphs in a single plot, we use `subplot`. From the plot we came to know,

Customer Care Calls: Slight positive skewed normal distribution with mode at 4.

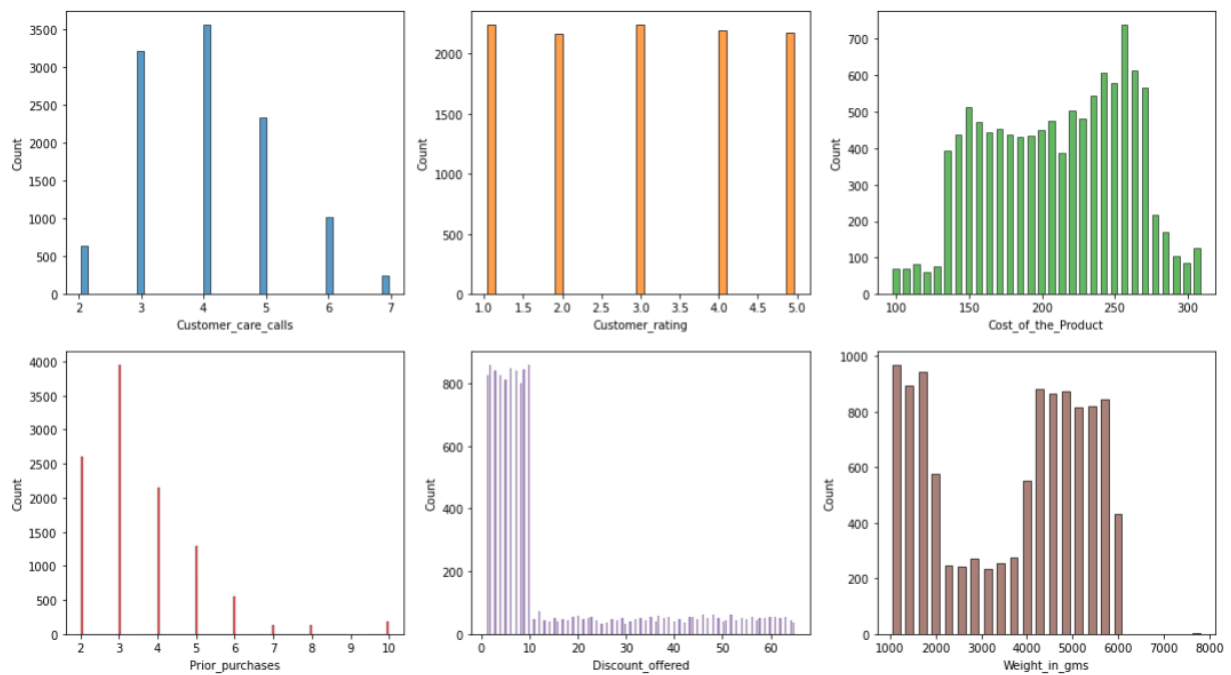
Customer Rating: Uniform distribution.

Cost of the product: 2 picks: smallest around 150, highest around 250.

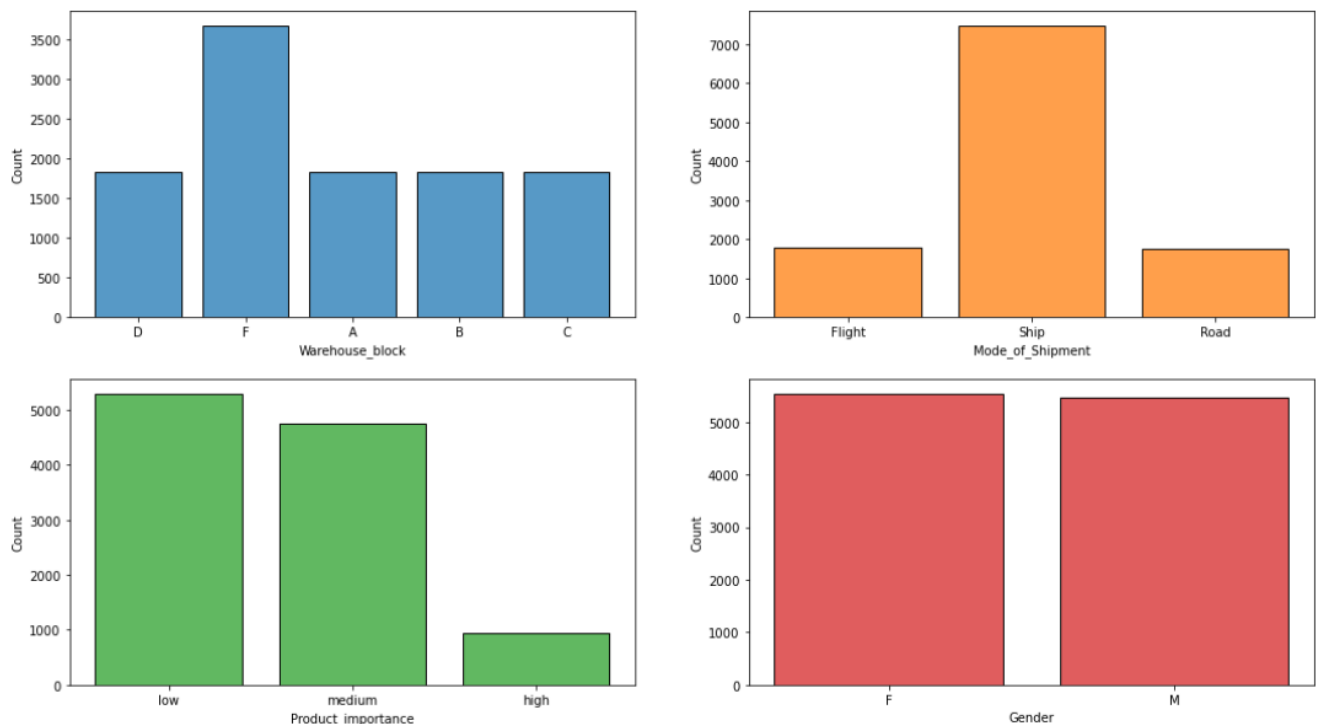
Prior Purchases: Positive skewed normal distribution, mode at 3.

Discount offered: Separated into 2 uniform distributions: 0 to 10 is predominant and then small amount from 10 to 65.

Weight: 3 zones: high from 1000 to 2000 and from 4000 to 6000. Low from 2000 to 4000.



In our dataset we have some categorical features. With the countplot function, we are going to count the unique category in those features. We have created a dummy data frame with categorical features. With for loop and subplot we have plotted this below graph.



Activity 2.2: Bivariate analysis

To find the relation between two features we use bivariate analysis. Here we are visualizing the relationship between drug & BP, drug & sex and drug & cholesterol.

- Countplot is used here. As a 1st parameter we are passing x value and as a 2nd parameter we are passing hue value.
- From the below plot you can understand that drugA and drugB is not preferred to low and normal BP patients. DrugC is preferred only to low BP patients.
- By third graph we can understand, drugC is not preferred to normal cholesterol patients.

Warehouse: Blocks A, B, C, D are equilibrated while block F is predominant (1/2 ratio).

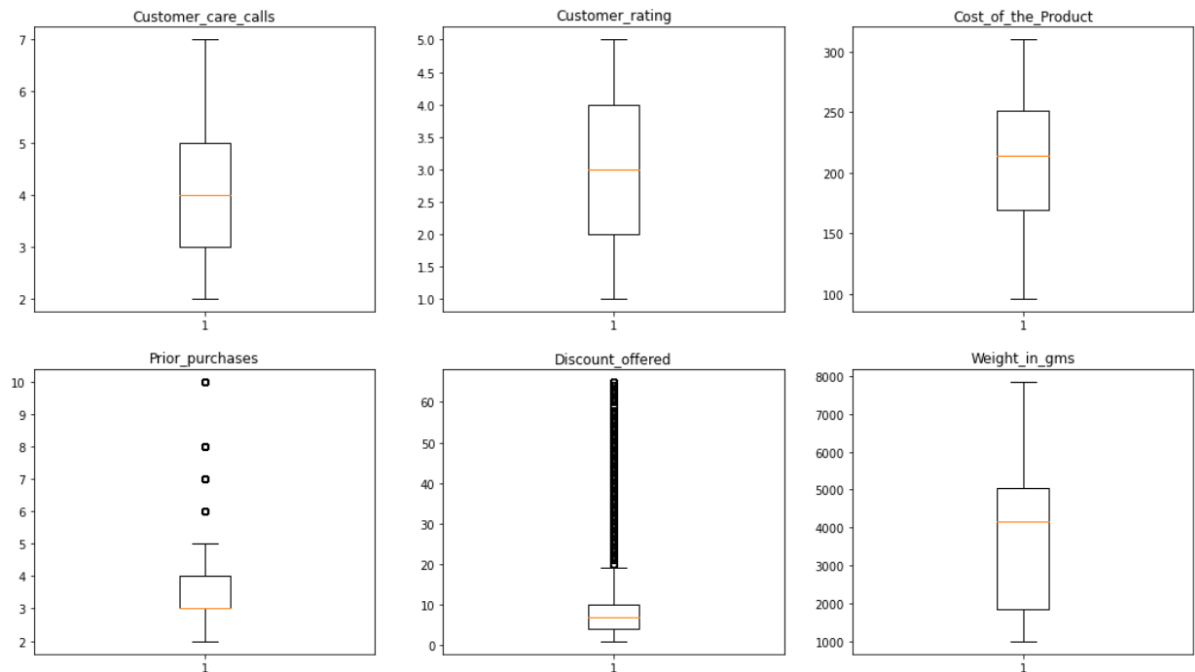
Shipment: Flight and Road have similar observations while Ship is predominant (1/4 ratio).

Importance: There is a majority of low and medium importances and a minority of high importances.

Genders: Both classes are balanced.

From the plots, we observe that very less number of orders are considered as important.

The most common mode of shipment is by Ship and products are packaged from warehouse F. Both genders order the products in a balanced way with orders by Females being slightly higher.



Outliers were found for 2 features (Discount_offered,Weight_in_gms) as visualized above using box plots. To be specific using IQR (inter quartile range) it was observed that,

Prior_purchases

1003 are over the upper bound: 5.5
0 are less than the lower bound: 1.5

Discount_offered

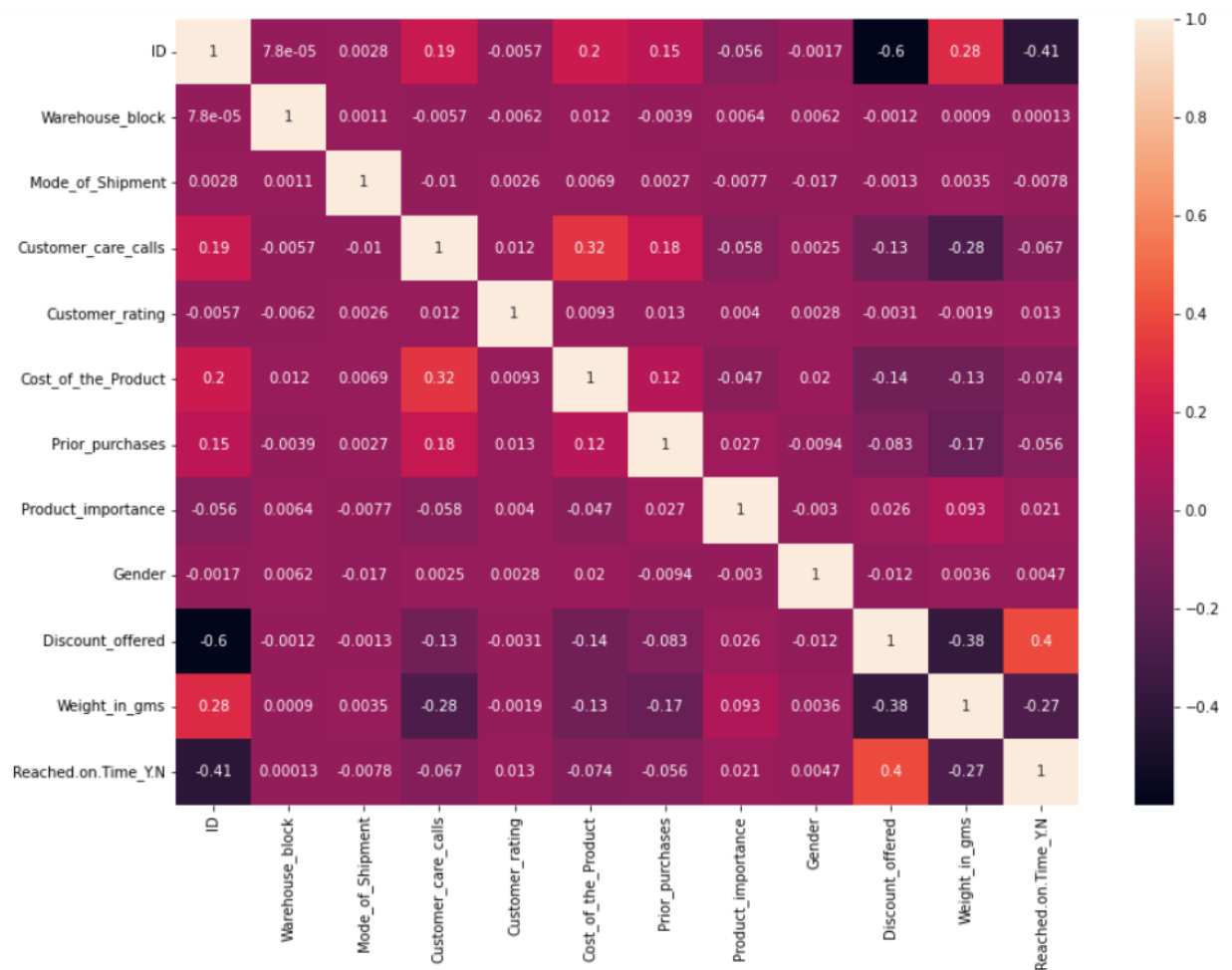
2262 are over the upper bound: 19.0
0 are less than the lower bound: -5.0

What this means is that there are 1003 and 2263 outliers for Prior_purchases and Discount_offered variables of the dataset.

Activity 2.3: Multivariate analysis

In simple words, multivariate analysis is to find the relation between multiple features. Here we have used heatmap from seaborn package.

- From the below image, we came to a conclusion that the product discount is the feature that most highly correlates to if a product is delivered on time and Number of calls and product cost are also highly correlated among other variables.



Splitting data into train and test

Now let's split the Dataset into train and test sets. First split the dataset into x and y and then split the data set.

Here x and y variables are created. On x variable, data is passed with dropping the target variable. And on y target variable is passed. For splitting training and testing data we are using `train_test_split()` function from `sklearn`. As parameters, we are passing x, y, `test_size`, `random_state` shuffle.

```
x_train,x_test,y_train,y_test=train_test_split(data.drop(columns=['ID', 'Reached.on.Time_Y.N']),data['Reached.on.Time_Y.N'])
print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)
```

Python

```
(8799, 10)
(2200, 10)
(8799,)
(2200,)
```


Milestone 4: Model Building

Activity 1: Training the model in multiple algorithms

Now our data is cleaned and it's time to build the model. We can train our data on different algorithms. For this project we are applying seven classification algorithms. The best model is saved based on its performance.

Activity 1.1: Writing function to train the models

A function named `models_eval_mm` is created and train, test data are passed as parameters. In the function, logistic regression, logistic regression cv, XGBclassifier, RidgeClassifier, KNN classifier, Random forest classifier and SVC algorithms are initialised and training data is passed to the model with `.fit()` function. Test data is predicted with `.predict()` function and saved in a new variable. For evaluating the model, train and test scores are used.

```
def models_eval_mm(x_train,y_train, x_test,y_test):
    lg=LogisticRegression (random_state=1234)
    lg.fit(x_train, y_train)
    print('--Logistic Regression')
    print('Train Score:', lg.score(x_train, y_train))
    print('Test Score:',lg.score(x_test,y_test))
    print()
    lcv= LogisticRegressionCV (random_state=1234)
    lcv.fit(x_train,y_train)
    print('--Logistic Regression CV')
    print('Train Score:',lcv.score(x_train,y_train))
    print()
    print('Test Score:',lcv.score(x_test,y_test))
    print('--XGBoost')
    xgb = XGBClassifier(random_state=1234)
    xgb.fit(x_train,y_train)
    print('Train Score:', xgb.score(x_train,y_train))
    print('Test Score:xgb',xgb.score(x_test,y_test))
    print()
    print('--Ridge Classifier')
    rg = RidgeClassifier(random_state=1234)
    rg.fit(x_train,y_train)
    print('Train Score:', rg.score(x_train, y_train))
    print('Test Score:',rg.score(x_test,y_test))
    print()
    print('--KNN')
    knn = KNeighborsClassifier()
    knn.fit(x_train,y_train)
    print('Train Score:',knn.score(x_train,y_train))
    print('Test Score:',knn.score(x_test,y_test))
    print()
    print('--Random Forest')
    rf = RandomForestClassifier (random_state=1234)
    rf.fit(x_train,y_train)
    print('Train Score:', rf.score(x_train,y_train))
    print('Test Score:',rf.score(x_test,y_test))
    print()
    print('--SVM classifier')
    svc = svm.SVC(random_state=1234)
    svc.fit(x_train,y_train)
    print('Train Score:', svc.score(x_train,y_train))
    print('Test Score:',svc.score(x_test,y_test))
    print()
    return lg,lcv, xgb, rg, knn, rf, svc
```

Python

Activity 1.2: Calling the function

The function is called by passing the train, test variables. The models are returned and stored in variables as shown below. Clearly, we can see that the models are not performing well on the data. So, we'll optimise the hyperparameters of models using `GridsearchCV`.

```
lg, lcv, xgb, rg, knn, rf, svc=models_eval_mm(x_train_normalized, y_train,x_test_normalized,y_test)

--Logistic Regression
Train Score: 0.5976815547221275
Test Score: 0.5927272727272728

--Logistic Regression CV
Train Score: 0.6422320718263439

Test Score: 0.6409090909090909
--XGBoost
Train Score: 0.9362427548585066
Test Score:xgb 0.6681818181818182

--Ridge Classifier
Train Score: 0.5976815547221275
Test Score: 0.5927272727272728

--KNN
Train Score: 0.7756563245823389
Test Score: 0.6336363636363637

--Random Forest
Train Score: 1.0
Test Score: 0.6686363636363636

--SVM classifier
Train Score: 0.5976815547221275
Test Score: 0.5927272727272728
```

Activity 2: Testing the model

Here we have tested with Random forest algorithm. You can test with all algorithm. With the help of predict() function.

```
lg.predict(x_test_normalized.iloc[0].values.reshape(1, -1))

c:\Users\Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\base.py:493: UserWarning: X
warnings.warn(
array([1], dtype=int64)
```

Milestone 5: Performance Testing & Hyperparameter Tuning

Activity 1: Testing model with multiple evaluation metrics

Multiple evaluation metrics means evaluating the model's performance on a test set using different performance measures. This can provide a more comprehensive understanding of the model's strengths and weaknesses. We are using evaluation metrics for classification tasks including accuracy, precision, recall, support and F1-score.

Activity 1.1: Compare the model

For comparing the above four models, the eval() function is defined.

```
Codumate: Options | Test this function
def eval(name, model):
    y_pred=model.predict(x_test_normalized)
    result =[]
    result.append(name)
    result.append("{:.2f}".format(accuracy_score(y_test, y_pred)*100))
    result.append("{:.2f}".format(f1_score(y_test, y_pred)*100))
    result.append("{:.2f}".format(recall_score(y_test, y_pred)*100))
    result.append("{:.2f}".format(precision_score(y_test, y_pred)*100))
    return result

model_list={'logistic regression': lg,
            'logistic regression CV':lcv,
            'XGBoost':xgb,
            'Ridge classifier':rg,
            'KNN':knn,
            "Random Forest":rf,
            "Support Vector Classifier":svc}

model_eval_info=[]
for i in model_list.keys():
    model_eval_info.append(eval(i,model_list[i]))
model_eval_info_df = pd.DataFrame(model_eval_info, columns=['Name', 'Accuracy', 'F1_Score', 'Recall', 'Precision'])

model_eval_info_df.to_csv("model_eval.csv", index=False)
from IPython.display import display, HTML

display(HTML(model_eval_info_df.to_html(index=False)))
```

[26]

...

Name	Accuracy	F1_Score	Recall	Precision
logistic regression	59.27	74.43	100.00	59.27
logistic regression CV	64.09	67.08	61.73	73.45
XGBoost	66.82	71.03	68.63	73.60
Ridge classifier	59.27	74.43	100.00	59.27
KNN	63.36	68.27	66.49	70.15
Random Forest	66.86	69.64	64.11	76.21
Support Vector Classifier	59.27	74.43	100.00	59.27

After calling the function, the results of models are displayed as output. From the trained models random forest is performing well

Activity 2: Comparing model accuracy before & after applying hyperparameter tuning (Hyperparameter tuning is optional. For this project it is not required.)

We will perform hyper parameter tuning on models that are better performing which are SVC, Random Forest, XGBoost and Logistic Regression CV. First we'll discover the existing hyperparameters of these models as below.

XGboost

```
xgb = XGBClassifier(learning_rate=0.5, n_estimators=100, objective='binary:logistic', nthread=3)

# Define the parameter grid for XGBoost
params = {
    'min_child_weight': [10, 20],
    'gamma': [1.5, 2.0, 2.5],
    'colsample_bytree': [0.6, 0.8, 0.9],
    'max_depth': [4, 5, 6]
}

# Initialize GridSearchCV for XGBoost
fitmodel = GridSearchCV(xgb, param_grid=params, cv=5, refit=True, scoring="accuracy", n_jobs=-1, verbose=3)

# Fit the model using the normalized training data
fitmodel.fit(x_train_normalized, y_train)

# Print the best estimator, parameters, and score
print("Best Estimator:", fitmodel.best_estimator_)
print("Best Parameters:", fitmodel.best_params_)
print("Best Score:", fitmodel.best_score_)

Fitting 5 folds for each of 54 candidates, totalling 270 fits
Best Estimator: XGBClassifier(base_score=None, booster=None, callbacks=None,
                               colsample_bylevel=None, colsample_bynode=None,
                               colsample_bytree=0.9, device=None, early_stopping_rounds=None,
                               enable_categorical=False, eval_metric=None, feature_types=None,
                               gamma=2.0, grow_policy=None, importance_type=None,
                               interaction_constraints=None, learning_rate=0.5, max_bin=None,
                               max_cat_threshold=None, max_cat_to_onehot=None,
                               max_delta_step=None, max_depth=4, max_leaves=None,
                               min_child_weight=20, missing=nan, monotone_constraints=None,
                               multi_strategy=None, n_estimators=100, n_jobs=None, nthread=3,
                               num_parallel_tree=None, ...)
Best Parameters: {'colsample_bytree': 0.9, 'gamma': 2.0, 'max_depth': 4, 'min_child_weight': 20}
Best Score: 0.6763268127551811
```

Random Forest

```
rf = RandomForestClassifier(random_state=1234)

rf_param_grid = {
    'n_estimators': [200, 300, 500],
    'criterion': ['entropy', 'gini'],
    'max_depth': [7, 8, 60, 80, 100],
    'max_features': ['auto', 'sqrt', 'log2']
}

rf_cv = GridSearchCV(rf, param_grid=rf_param_grid, cv=7, scoring="accuracy", n_jobs=-1, verbose=3)

rf_cv.fit(x_train_normalized, y_train)

print("Best Score:", rf_cv.best_score_)
print("Best Parameters:", rf_cv.best_params_)
```

Python

Fitting 7 folds for each of 90 candidates, totalling 630 fits
[C:\Users\Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\model_selection_validation.py:549: FitFailedWarning:](#)
210 fits failed out of a total of 630.
The score on these train-test partitions for these parameters will be set to nan.
If these failures are not expected, you can try to debug them by setting error_score='raise'.

Below are more details about the failures:

133 fits failed with the following error:

Traceback (most recent call last):
File "C:\Users\Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\model_selection_validation.py", line 888, in _fit_and_score
estimator.fit(X_train, y_train, **fit_params)
File "C:\Users\Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\base.py", line 1466, in wrapper
estimator._validate_params()
File "C:\Users\Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\base.py", line 666, in _validate_params
validate_parameter_constraints(
File "C:\Users\Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\utils_param_validation.py", line 95, in validate_parameter_constraints
raise InvalidParameterError(
sklearn.utils._param_validation.InvalidParameterError: The 'max_features' parameter of RandomForestClassifier must be an int in the range [1, inf), a float in the range (0.0, 1.0], a str among ('log2', 'sqrt') or None. Got 'auto'

77 fits failed with the following error:

Traceback (most recent call last):
File "C:\Users\Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\model_selection_validation.py", line 888, in _fit_and_score
estimator.fit(X_train, y_train, **fit_params)
File "C:\Users\Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\base.py", line 1466, in wrapper
...
nan nan nan 0.65416525 0.65518889 0.6533697
0.65416525 0.65518889 0.6533697 nan nan nan
0.65416525 0.65518889 0.6533697 0.65416525 0.65518889 0.6533697]
warnings.warn(
Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#).
Best Score: 0.6801909307875896
Best Parameters: {'criterion': 'entropy', 'max_depth': 8, 'max_features': 'sqrt', 'n_estimators': 200}

Logistic Regression CV

```
Logistic Regression

lg = LogisticRegressionCV(n_jobs=-1, random_state=1234)

lg_param_grid = {
    'Cs': [6, 8, 10, 15, 20],
    'max_iter': [60, 80, 100]
}

lg_cv = GridSearchCV(lg, param_grid=lg_param_grid, cv=5, scoring='accuracy', n_jobs=-1, verbose=3)

lg_cv.fit(x_train_normalized, y_train)

print("Best Score:", lg_cv.best_score_)
print("Best Parameters:", lg_cv.best_params_)

Fitting 5 folds for each of 15 candidates, totalling 75 fits
Best Score: 0.6412080126853026
Best Parameters: {'Cs': 20, 'max_iter': 60}
```

To perform hyper parameter tuning use the following code as shown below.

```
HyperParameter Optimization
SVM

Cadmate: Options [Test this function]
import pandas as pd
from sklearn import svm
from sklearn.model_selection import GridSearchCV

# Define the SVM model
svc = svm.SVC(random_state=1234)

# Define the parameter grid
params = {
    'kernel': ['poly', 'rbf'],
    'C': [10, 15],
    'gamma': [4, 5],
    'tol': [1e-1, 1e-2, 1e-3]
}

# Initialize GridSearchCV
svc_cv = GridSearchCV(svc, param_grid=params, cv=5, refit=True, scoring="accuracy", n_jobs=-1, verbose=3)

# Fit the model
svc_cv.fit(x_train_normalized, y_train)

# Print the best estimator, parameters, and score
print("Best Estimator:", svc_cv.best_estimator_)
print("Best Parameters:", svc_cv.best_params_)
print("Best Score:", svc_cv.best_score_)

Fitting 5 folds for each of 24 candidates, totalling 120 fits
Best Estimator: SVC(C=15, gamma=5, random_state=1234, tol=0.01)
Best Parameters: {'C': 15, 'gamma': 5, 'kernel': 'rbf', 'tol': 0.01}
Best Score: 0.6676890278567368
```

```
XGB

xgb = XGBClassifier(learning_rate=0.5, n_estimators=100, objective='binary:logistic', nthread=3)

# Define the parameter grid for XGBoost
params = {
    'min_child_weight': [10, 20],
    'gamma': [1.5, 2.0, 2.5],
    'colsample_bytree': [0.6, 0.8, 0.9],
    'max_depth': [4, 5, 6]
}

# Initialize GridSearchCV for XGBoost
fitmodel = GridSearchCV(xgb, param_grid=params, cv=5, refit=True, scoring="accuracy", n_jobs=-1, verbose=3)

# Fit the model using the normalized training data
fitmodel.fit(x_train_normalized, y_train)

# Print the best estimator, parameters, and score
print("Best Estimator:", fitmodel.best_estimator_)
print("Best Parameters:", fitmodel.best_params_)
print("Best Score:", fitmodel.best_score_)

Fitting 5 folds for each of 54 candidates, totalling 270 fits
Best Estimator: XGBClassifier(base_score=None, booster=None, callbacks=None,
    colsample_bylevel=None, colsample_bynode=None,
    colsample_bytree=0.9, device=None, early_stopping_rounds=None,
    enable_categorical=False, eval_metric=None, feature_types=None,
    gamma=2.0, grow_policy=None, importance_type=None,
    interaction_constraints=None, learning_rate=0.5, max_bin=None,
    max_cat_threshold=None, max_cat_to_onehot=None,
    max_delta_step=None, max_depth=4, max_leaves=None,
    min_child_weight=20, missing=nan, monotone_constraints=None,
    multi_strategy=None, n_estimators=100, n_jobs=None, nthread=3,
    num_parallel_tree=None, ...)
Best Parameters: {'colsample_bytree': 0.9, 'gamma': 2.0, 'max_depth': 4, 'min_child_weight': 20}
Best Score: 0.6763268127551811
```

Logistic Regression

```
lg = LogisticRegressionCV(n_jobs=-1, random_state=1234)

lg_param_grid = {
    'Cs': [6, 8, 10, 15, 20],
    'max_iter': [60, 80, 100]
}

lg_cv = GridSearchCV(lg, param_grid=lg_param_grid, cv=5, scoring="accuracy", n_jobs=-1, verbose=3)

lg_cv.fit(x_train_normalized, y_train)

print("Best Score:", lg_cv.best_score_)
print("Best Parameters:", lg_cv.best_params_)
```

```
29)
... Fitting 5 folds for each of 15 candidates, totalling 75 fits
Best Score: 0.6412089126053025
Best Parameters: {'Cs': 20, 'max_iter': 60}
```

Random Forest

```
rf = RandomForestClassifier(random_state=1234)

rf_param_grid = {
    'n_estimators': [200, 300, 500],
    'criterion': ['entropy', 'gini'],
    'max_depth': [7, 9, 60, 80, 100],
    'max_features': ['auto', 'sqrt', 'log2']
}

rf_cv = GridSearchCV(rf, param_grid=rf_param_grid, cv=7, scoring="accuracy", n_jobs=-1, verbose=3)

rf_cv.fit(x_train_normalized, y_train)

print("Best Score:", rf_cv.best_score_)
print("Best Parameters:", rf_cv.best_params_)
```

Fitting 7 folds for each of 90 candidates, totalling 630 fits

210 fits failed out of a total of 630.
The score on these train-test partitions for these parameters will be set to nan.
If these failures are not expected, you can try to debug them by setting `error_score='raise'`.

Below are more details about the failures:

130 fits failed with the following error:

Traceback (most recent call last):
File "c:\Users\Arjun\AppData\Local\Programs\Python\Python311\lib\site-packages\sklearn\model_selection_validation.py", line 888, in _fit_and_score
estimator.fit(X_train, y_train, **fit_params)
File "c:\Users\Arjun\AppData\Local\Programs\Python\Python311\lib\site-packages\sklearn\base.py", line 1466, in wrapper
estimator._validate_params()
File "c:\Users\Arjun\AppData\Local\Programs\Python\Python311\lib\site-packages\sklearn\base.py", line 666, in _validate_params
validate_parameter_constraints(
File "c:\Users\Arjun\AppData\Local\Programs\Python\Python311\lib\site-packages\sklearn\utils_param_validation.py", line 95, in validate_parameter_constraints
raise InvalidParameterError(
sklearn.utils._param_validation.InvalidParameterError: The 'max_features' parameter of RandomForestClassifier must be an int in the range [1, inf), a float in the range (0.0, 1.0], a str among ('log2', 'sqrt') or None. Got 'auto'

77 fits failed with the following error:

Traceback (most recent call last):
File "c:\Users\Arjun\AppData\Local\Programs\Python\Python311\lib\site-packages\sklearn\model_selection_validation.py", line 888, in _fit_and_score
estimator.fit(X_train, y_train, **fit_params)
File "c:\Users\Arjun\AppData\Local\Programs\Python\Python311\lib\site-packages\sklearn\base.py", line 1466, in wrapper
...
nan nan nan 0.65416525 0.65518809 0.6533607
0.65416525 0.65518809 0.6533607 nan nan nan
0.65416525 0.65518809 0.6533607 0.65416525 0.65518809 0.6533607]
warnings.warn(
Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output settings.
Best Score: 0.6081909307075890
Best Parameters: {'criterion': 'entropy', 'max_depth': 0, 'max_features': 'sqrt', 'n_estimators': 200}

Model evaluation before hyper parameter tuning

	Name	Accuracy	F1_Score	Recall	Precision
	logistic regression	59.27	74.43	100.00	59.27
	logistic regression CV	64.09	67.08	61.73	73.45
	XGBoost	66.82	71.03	68.63	73.60
	Ridge classifier	59.27	74.43	100.00	59.27
	KNN	63.36	68.27	66.49	70.15
	Random Forest	66.86	69.64	64.11	76.21
	Support Vector Classifier	59.27	74.43	100.00	59.27

Model evaluation after hyper parameter tuning

Model Evaluation Results:

	Name	Accuracy	F1_Score	Recall	Precision
	Logistic Regression	63.95	66.97	61.66	73.29
	Logistic Regression CV	64.09	67.08	61.73	73.45
	XGBoost	67.36	68.40	59.59	80.27
	Ridge Classifier	65.05	71.84	75.23	68.75
	KNN	64.55	68.82	66.03	71.87
	Random Forest	67.45	66.13	53.60	86.30
	Support Vector Classifier	66.77	68.29	60.35	78.62

Milestone 6: Model Deployment

Activity 1: Save the best model

Saving the best model and normalizer after comparing its performance using different evaluation metrics means selecting the model with the highest performance and saving its weights and configuration. This can be useful in avoiding the need to retrain the model every time it is needed and also to be able to use it in the future.

```
import pickle
```

```
filename = 'finalized_knn_regression.pkl'  
pickle.dump(KNN_model, open(filename, 'wb'))
```

Activity 2: Integrate with Web Framework

In this section, we will be building a web application that is integrated to the model we built. An UI is provided for the uses where he has to enter the values for predictions. The enter values are given to the saved model and prediction is showcased on the UI.

This section has the following tasks

- Building HTML Pages
- Building server-side script
- Run the web application

Activity 2.1: Building Html Pages:

For this project create two HTML files namely

- index.html
- inner-page.html
- portfolio-details.html

and save them in the templates folder.

It is not necessary to follow the exact format as above so feel free to use whatever templates or format you like. Be creative!

Activity 2.2: Build Python code:

Import the libraries

```
from flask import Flask, render_template, request  
import numpy as np  
import pickle
```

Load the saved model. Importing the flask module in the project is mandatory. An object of Flask class is our WSGI application. Flask constructor takes the name of the current module (`__name__`) as argument.

```
app = Flask(__name__)  
  
model = pickle.load(open('finalized_knn_regression.pkl', 'rb'))
```

Render HTML page:

```
@app.route("/")  
def home():  
    return render_template("home.html")
```

Here we will be using a declared constructor to route to the HTML page which we have created earlier.

In the above example, '/' URL is bound with the index.html function. Hence, when the home page of the web server is opened in the browser, the html page will be rendered. Whenever you enter the values from the html page the values can be retrieved using POST Method.

Retrieves the value from UI:


```

from flask import Flask, render_template, request
import numpy as np
import pickle

app = Flask(__name__)

# Load the trained model
model = pickle.load(open('finalized_knn_regression.pkl', 'rb'))

Codiumate: Options | Test this function
@app.route("/")
def home():
    return render_template("home.html")

Codiumate: Options | Test this function
@app.route("/result", methods=["POST"])
def submit():
    # Extract form data
    global cost_of_the_product, weight_in_gms, discount_offered
    global warehouse_block_A, warehouse_block_B, warehouse_block_C, warehouse_block_D, warehouse_block_F
    global mode_of_shipment_Flight, mode_of_shipment_Road, mode_of_shipment_Ship
    global product_importance_high, product_importance_low, product_importance_medium
    global customer_care_calls, prior_purchases, gender

    if request.method == "POST":
        warehouse_block = request.form["warehouse_block"]
        if warehouse_block == "A":
            warehouse_block_A = 1
            warehouse_block_B = 0
            warehouse_block_C = 0
            warehouse_block_D = 0
            warehouse_block_F = 0
        elif warehouse_block == "B":
            warehouse_block_A = 0
            warehouse_block_B = 1
            warehouse_block_C = 0
            warehouse_block_D = 0
            warehouse_block_F = 0
        elif warehouse_block == "C":
            warehouse_block_A = 0
            warehouse_block_B = 0
            warehouse_block_C = 1
            warehouse_block_D = 0
            warehouse_block_F = 0
        elif warehouse_block == "D":
            warehouse_block_A = 0
            warehouse_block_B = 0
            warehouse_block_C = 0
            warehouse_block_D = 1
            warehouse_block_F = 0
        elif warehouse_block == "F":
            warehouse_block_A = 0
            warehouse_block_B = 0
            warehouse_block_C = 0
            warehouse_block_D = 0
            warehouse_block_F = 1

```

```

mode_of_shipment = request.form["mode_of_shipment"]
if mode_of_shipment == "Ship":
    mode_of_shipment_Flight = 0
    mode_of_shipment_Road = 0
    mode_of_shipment_Ship = 1
elif mode_of_shipment == "Flight":
    mode_of_shipment_Flight = 1
    mode_of_shipment_Road = 0
    mode_of_shipment_Ship = 0
elif mode_of_shipment == "Road":
    mode_of_shipment_Flight = 0
    mode_of_shipment_Road = 1
    mode_of_shipment_Ship = 0

product_importance = request.form["product_importance"]
if product_importance == "low":
    product_importance_high = 0
    product_importance_low = 1
    product_importance_medium = 0
elif product_importance == "medium":
    product_importance_high = 0
    product_importance_low = 0
    product_importance_medium = 1
elif product_importance == "high":
    product_importance_high = 1
    product_importance_low = 0
    product_importance_medium = 0

customer_care_calls = int(request.form["customer_care_calls"])
gender = int(request.form["gender"])
prior_purchases = int(request.form["prior_purchases"])
cost_of_the_product = int(request.form["cost_of_the_product"])
discount_offered = int(request.form["discount_offered"])
weight_in_gms = int(request.form["weight_in_gms"])

# Prepare the feature array for prediction
x = np.array([cost_of_the_product, weight_in_gms, discount_offered,
              warehouse_block_A, warehouse_block_B, warehouse_block_C, warehouse_block_D, warehouse_block_F,
              mode_of_shipment_Flight, mode_of_shipment_Road, mode_of_shipment_Ship,
              product_importance_high, product_importance_low, product_importance_medium,
              customer_care_calls, prior_purchases, gender])
x = x.reshape((1, -1))

# Make the prediction (regression output)
prediction = model.predict(x)[0]
prediction = int(prediction*100)
# Render result.html with the prediction
return render_template('home.html', prediction=prediction)

if __name__ == "__main__":
    app.run(debug=True)

```

Here we are routing our app to predict() function. This function retrieves all the values from the HTML page using Post request. That is stored in an array. This array is passed to the model.predict() function. This function returns the prediction. And this prediction value will be rendered to the text that we have mentioned in the submit.html page earlier.

Main Function:

```
if __name__ == "__main__":  
    app.run(debug=True)
```

Activity 2.3: Run the web application

- Open anaconda prompt from the start menu
- Navigate to the folder where your python script is.
- Now type “python app.py” command
- Navigate to the localhost where you can view your web page.
- Click on the predict button from the top left corner, enter the inputs, click on the submit button, and see the result/prediction on the web.

```
Python 3.11.7 | packaged by Anaconda, Inc. | (main, Dec 15 2023, 18:05:47) [MSC v.1916 64 bit (AMD64)]  
Type "copyright", "credits" or "license" for more information.  
  
IPython 8.20.0 -- An enhanced Interactive Python.  
  
In [1]: runfile('C:/Users/Arjun/OneDrive/Desktop/VIT/Ecom(ayush)/intenshipProject/app.py', wdir='C:/  
Users/Arjun/OneDrive/Desktop/VIT/Ecom(ayush)/intenshipProject')  
* Serving Flask app 'app'  
* Debug mode: on  
C:\Users\Arjun\anaconda3\Lib\site-packages\sklearn\base.py:376: InconsistentVersionWarning: Trying to  
unpickle estimator KNeighborsRegressor from version 1.5.0 when using version 1.5.1. This might lead to  
breaking code or invalid results. Use at your own risk. For more info please refer to:  
https://scikit-learn.org/stable/model\_persistence.html#security-maintainability-limitations  
warnings.warn(  
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI  
server instead.  
* Running on http://127.0.0.1:5000  
Press CTRL+C to quit
```

Now, Go the web browser and write the localhost url (<http://127.0.0.1:4000>) to get the below result

Upon entering the data in input fields and clicking on submit we'll get the prediction as shown below.



On-Time Delivery Prediction

Let's check whether the package will reach on time.

Q1. Select the gender of the customer.

Q2. Select the warehouse block that is in charge.

Q3. Select the mode of shipment.

Q4. Select the product importance.

Q5. Select the number of calls that the customer has made.

Q6. Select the purchases that the customer had made.

Q7. Did the shipment arrive on time?

Q8. Select the cost of the product (in USD).

Q9. Select the discount given to the customer.

Q10. Select the weight (in grams) of the product.

Submit



On-Time Delivery Prediction

Let's check whether the package will reach on time.

Q1. Select the gender of the customer.

Q2. Select the warehouse block that is in charge.

Q3. Select the mode of shipment.

Q4. Select the product importance.

Q5. Select the number of calls that the customer has made.

Q6. Select the purchases that the customer had made.

Q7. Did the shipment arrive on time?

Q8. Select the cost of the product (in USD).

Q9. Select the discount given to the customer.

Q10. Select the weight (in grams) of the product.

Submit

Probability of Arriving on Time: 36%

