

# EXPERIMENT REPORT

<b>Name</b>	Monali Patil
<b>Project Name</b>	Customer-Vehicle-Repurchase-Prediction
<b>Deliverables</b>	<p>Model name: Logistic Regression Classifier</p> <p>Notebook name: MachineLearning_LogisticRegression_Exp2.ipynb</p> <p>Project Repo: <a href="https://github.com/monalipatil/MachineLearning-Customer-Vehicle-Repurchase-Prediction.git">https://github.com/monalipatil/MachineLearning-Customer-Vehicle-Repurchase-Prediction.git</a></p>

## 1. EXPERIMENT BACKGROUND

Provide information about the problem/project such as the scope, the overall objective, expectations. Lay down the goal of this experiment and what are the insights, answers you want to gain or level of performance you are expecting to reach.

### 1.a. Business Objective

Explain clearly what is the goal of this project for the business. How will the results be used? What will be the impact of accurate or incorrect results?

Objective: The aim of this project is to develop a binary classifier model that can predict customers who are likely to purchase a new vehicle. To achieve this, a dataset of car repurchases will be used for this binary classification problem and through data exploration, cleaning, feature scaling, selection of appropriate algorithm, hyperparameters and employing various techniques such as cross-validation, hyperparameter tuning, feature importance etc., to improve the accuracy of the model.

Application: The outcomes of this model can be used to identify the customers who are most likely to buy a new vehicle, and this information can be used to focus marketing efforts on those potential customers which can lead to an effective and efficient promotional campaign.

Impact: The model's accurate results can help the company to target the right customers who are more likely to buy the vehicle through a marketing campaign which could save marketing costs and the campaign could generate higher sales with fewer marketing expenses. However, incorrect results of the model can result in wasted marketing efforts and resources and thus, the marketing campaign may fail to achieve the desired results and sales figures.

<b>1.b. Hypothesis</b>	<p>Present the hypothesis you want to test, the question you want to answer or the insight you are seeking. Explain the reasons why you think it is worthwhile considering it,</p> <p>Hypothesis: The objective of this experiment is to effectively identify if present customers are likely to buy a new vehicle, and to carry out this experiment, the Logistic Regression Classifier algorithm has been employed for conducting and evaluating the performance to determine if this algorithm is suitable for this Binary Classification task.</p> <p>Rationale: Considering the business objective of determining potential customers leads for a marketing campaign, there is no good theory to map and select a suitable algorithm for this binary classification problem, therefore performing different experiments to discover which algorithm and algorithm configuration results in the best performance for this binary classification task.</p>
<b>1.c. Experiment Objective</b>	<p>Detail what will be the expected outcome of the experiment. If possible, estimate the goal you are expecting. List the possible scenarios resulting from this experiment.</p> <p>The expected outcome of this experiment would be the development of a binary classifier model utilizing Logistic Regression Classifier algorithm that can accurately predict which customers are more likely to buy a new vehicle. The model would be trained on a dataset of current customer information and their buying behavior, and it would use Logistic Regression Classifier algorithm to identify patterns and make predictions about future buying behaviour of the customers.</p> <p>The goal of this experiment would be to identify potential customers for a promotional marketing campaign aimed at increasing vehicle sales, employing Logistic Regression Classifier algorithm. By focusing on customers who are more likely to buy a new vehicle, the marketing campaign could be more targeted and effective.</p> <p>The possible scenarios resulting from this experiment include following:</p> <ul style="list-style-type: none"> <li>• The model accurately predicts which customers are likely to buy a new vehicle, allowing the marketing team to focus on these potential customers and increase sales.</li> <li>• The model has a high false positive rate, indicating that the model predicts some customers will purchase a new vehicle when they actually won't purchase. This could result in wasted marketing resources and a lower return on investment.</li> <li>• The model has a high false negative rate, indicating that the model fails to predict some customers will purchase a new vehicle when they actually will. This could result in missed opportunities for sales and revenue.</li> <li>• The model is not accurate enough to be useful for predicting customer behavior, and the experiment is unsuccessful. In this case, the marketing team may need to explore other methods for identifying potential customers for their promotional campaign.</li> </ul>

## 2. EXPERIMENT DETAILS

Elaborate on the approach taken for this experiment. List the different steps/techniques used and explain the rationale for choosing them.

### 2.a. Data Preparation

Describe the steps taken for preparing the data (if any). Explain the rationale why you had to perform these steps. List also the steps you decided to not execute and the reasoning behind it. Highlight any step that may potentially be important for future experiments.

To prepare for using the binary classification algorithm, the data for current customers was explored, cleaned up and prepared through the subsequent activities.

- Data Understanding

**1] Loading Data:** To use and create a binary classification model, imported data from a CSV file into the pandas dataframe.

**2] Exploring Data:** Explored and analysed current customer's data utilizing various following pandas functionalities to understand and identify patterns to further determine potential customers who are more likely to purchase a new vehicle.

\* Additionally, to ensure the quality of the data to be utilized by the model analyzed.

- Missing/null values.
- Duplicate records.
- Outliers for numerical features.
- Data distribution for various categorical and numerical features.
- Distinct values of the categorical features.

a) df.head(): Examining initial observations of the dataset.

b) df.shape(): Examining the dimension of the dataset.

c) df.columns: Inspecting features name.

d) df.info(): Inspecting the summary information of the attributes of the dataset.

e) df.describe(): Examining the statistical information of integer variables of the dataset.

f) df.describe(include='all'): Examining statistical summary information for all variables of the dataset across different data types.

g) df.isnull().sum(): Examining whether there are any null values in the dataset.

h) df.duplicated().sum(): Examining whether there are any missing values in the dataset.

- Data Preparation

### 3] Treating Missing Values

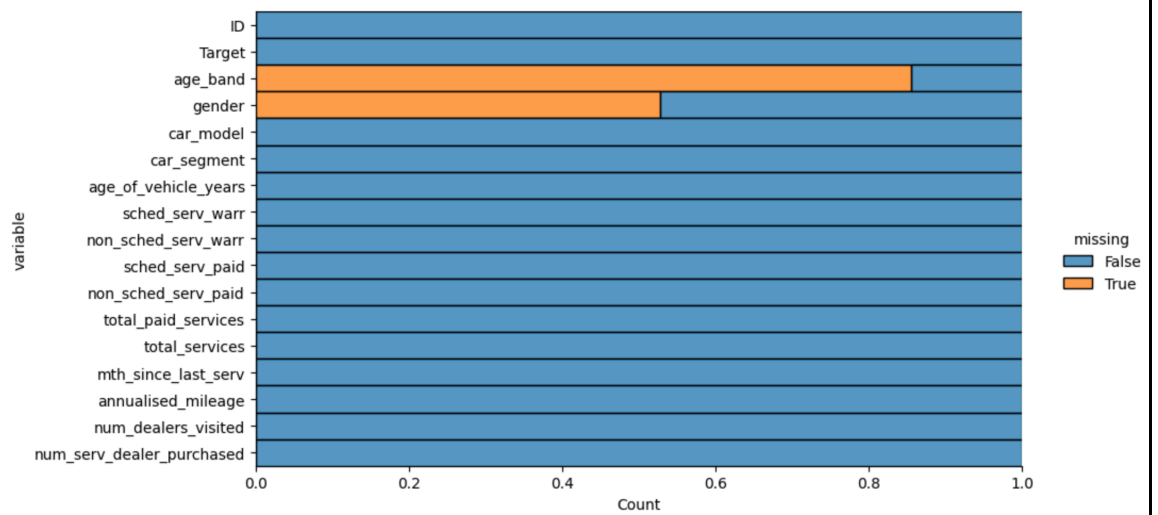


Figure 1: Missing values present in all attributes of the dataset.

\* The chart from Figure 1, indicates that only two variables, 'age\_band' with around 85% and 'gender' with over 50% contain missing values. However, there are no missing values for the rest of the features.

\* As machine learning algorithms are incapable of processing features with missing values, we have two options to address these missing values.

i) Replace missing values with the relevant value: Using the most frequent occurring (mode) value from the feature to replace missing values in categorical variables may not accurately represent the distribution of data in those variables 'age\_band' and 'gender', especially when large portion of the data is missing. This approach could also create an imbalance in the data by introducing a high number of imputed mode values, potentially leading to biased results.

ii) Eliminate the columns containing missing values entirely: It would be suitable to remove these two 'age\_band' and 'gender' features since replacing any value would significantly distort the data due to the high percentage of missing values in these variables.

\* Therefore, removed 'age\_band' and 'gender' features from the vehicle repurchase dataset using pandas drop() function and specifying these two feature names.

#### 4] Transforming Categorical Data into Numerical

Attribute Name	Distinct/Unique Values	Count of Distinct Values
car_model	model_1, model_2, model_3, model_5, model_6, model_4, model_7, model_8, model_9, model_10, model_11, model_13, model_12, model_14, model_15, model_16, model_17, model_18, model_19	19
car_segment	LCV, Small/Medium, Large/SUV, Other	4

Table 1: Distinct values present in 'car\_model' and 'car\_segment' attributes of the dataset.

\* The table 1 illustrates that the features 'car\_model' and 'car\_segment' contain categorical data, with 19 distinct car models and 4 different vehicle segments respectively. However, this data

cannot be used directly in binary classification algorithms for machine learning because these algorithms operate on mathematical equations that require numerical inputs.

\* Additionally, these features 'car\_model' and 'car\_segment' hold significant information regarding the model and category of vehicles that customers have purchased. This information can be used to predict whether the customer is inclined to purchase a new vehicle.

#### One-Hot Encoding

\* Therefore, it is necessary to convert the categorical data of these features into numerical data. To achieve this, employing one-hot encoding method, which creates binary columns for every distinct category in the feature.

\* This method is preferred as there is no inherent order among the category values of the feature, and do not want to create any random relationships between them and prevent model from assuming natural ordering among these distinct categories that may suffer from model bias.

### **5] Selecting Target Variable & Features**

\* The 'TARGET' feature is the target variable representing customers who have bought more than one vehicle with class 1 and those who have only purchased one vehicle with class 0. This, along with the predictor features, are separated into 'X' and 'y' variables to build a binary classification model.

### **6] Splitting Data into Different Sets: Training, Validation and Testing**

\* Creating a validation set provides adaptability to carry out multiple experiments and allows the comparison of the model's performance on the training set and the validation set to determine if the model is overfitting or underfitting.

\* Furthermore, using a validation set, we can adjust model's hyperparameters and retrain the model until it performs well on the validation set. After we have optimized the model's performance on the validation set, we can then apply the model on the test set to assess its ability to generalize to new and unseen data.

\* Therefore, utilizing train\_test\_split() function from sklearn.model\_selection module of sklearn library, the car repurchase dataset is split into training (60%), validation (20%), and testing (20%) sets to leverage the flexibility of multiple experimentation.

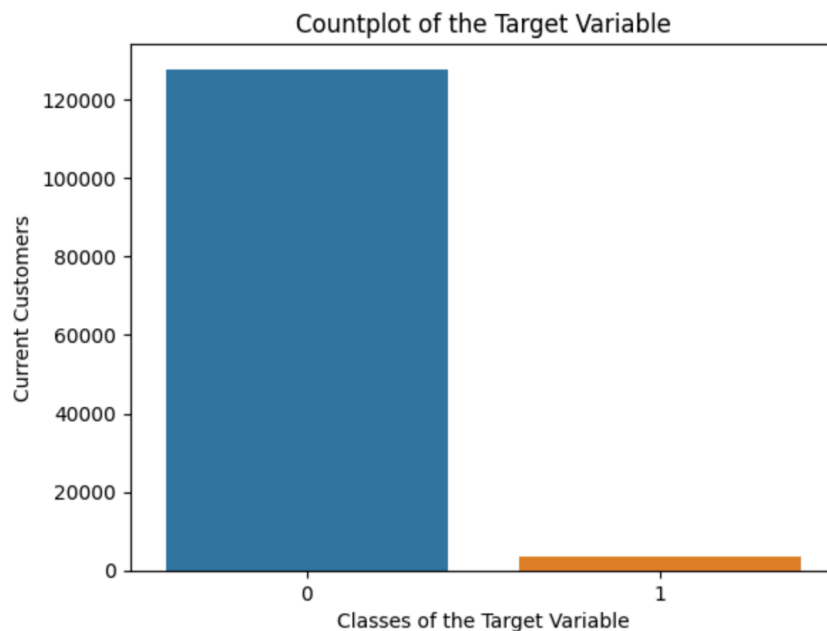


Figure 2: Examining the quantity of each class of the target variable from the actual dataset.

\* Additionally, for any (binary) classification problem, it is essential to examine the frequencies/rates of every class within the target variable and ensure that the splitting of data into different sets matches as similar to the actual data spread.

\* As indicated in the Figure 2, the car repurchase dataset is highly imbalanced, it is necessary to obtain as similar distribution of distinct classes of the target variable for all of the different sets as original dataset. After splitting process, following are the frequencies of the target variable classes.

Dataset/Frequency Rates	Class 1	Class 0
Actual Data:	0.97	0.026
Training Set:	0.97	0.026
Validation Set:	0.97	0.026
Testing Set:	0.97	0.028

## 7] Features Scaling

\* Machine learning algorithms often use some form of distance calculation to determine the relationship between different features in the dataset. If the features are on different scales, the algorithm may prioritize some high value range features over others which may have more significant information, leading to biased results.

\* In addition, the values for the numerical features are divided into ten deciles, ranging from 1 to 10, while categorical data that has been converted into numerical features, consist of either 0 or 1.

\* Therefore, performing scaling for all the features so that all the features values are at same level and that the algorithm can use all the information from the dataset features to learn generalised patterns, identify buying behaviour and make accurate predictions.

#### StandardScaler

\* Employing StandardScaler class which is imported from sklearn.preprocessing module to scale and bring all features values at same level. Choosing this method because it does not change the shape of the data distribution and preserves outliers as it scales the data based on the mean of 0 and standard deviation of 1, of the entire dataset, rather than individual datapoints.

The following are some crucial steps that could be important for any of the classification experiments in future.

- It is important to thoroughly examine and handle missing values in order to avoid introducing biases in the model.
- Outliers and duplicate values should be carefully analysed and reviewed from business perspectives and treated accordingly.
- The decision to use LabelEncoder, OrdinalEncoder or OneHotEncoder for converting categorical data into numerical form should depend on various factors, including the nature of the categorical values and whether or not there is any inherent order among them. Careful consideration is necessary to determine the appropriate method.
- Similarly, the selection of feature scaling method - whether to use MinMaxScaler, MaxAbsScaler, or StandardScaler - should be made based on whether or not there is an inherent order in the data, as well as other factors such as the need to maintain the integrity of outliers.
- For classification problems and especially for imbalanced data, is essential to ensure the frequency of each class in the target variable in each set is representative of the actual data. The "stratify=y" parameter could be used during the splitting process if the class distribution is not similar to the actual data. This will help to maintain the integrity of the data and ensure accurate results.
- To prevent the model from overfitting on specific data points and to enable it to learn generalized patterns, it is important to eliminate any identifier attributes.

## 2.b. Feature Engineering

Describe the steps taken for generating features (if any). Explain the rationale why you had to perform these steps. List also the feature you decided to remove and the reasoning behind it. Highlight any feature that may potentially be important for future experiments

	ID
count	131337.000000
mean	77097.384180
std	44501.636704
min	1.000000
25%	38563.000000
50%	77132.000000
75%	115668.000000
max	154139.000000

Figure 3: Statistical summary of 'ID' attribute.

\* The 'ID' attribute seems to be an identifier and its values ranges from 1 to 154139 as indicated in Figure 3. It does not have any predictive value as it contains a unique value for each observation, which means that its values do not contribute to any patterns or trends in the data.

\* And including them in the analysis can lead to overfitting straight a way, where model would fit these specific values from the feature rather than the underlying generic patterns in the data.

\* Therefore, removed 'ID' attribute from the vehicle repurchase dataset using pandas drop() function and specifying this feature name.

Note: Removal of 'age\_band' and 'gender' features is stated in earlier part 3] Treating Missing Values of 2.a Data Preparation section.

## 2.c. Modelling

Describe the model(s) trained for this experiment and why you choose them. List the hyperparameter tuned and the values tested and also the rationale why you choose them. List also the models you decided to not train and the reasoning behind it. Highlight any model or hyperparameter that may potentially be important for future experiments

The model trained for this experiment is one of the Binary Classification algorithms, Logistic Regression Classifier model to accurately determine whether current customers are probable to purchase a new vehicle. The model is trained, and the performance is evaluated to determine if this algorithm is suitable for the Binary Classification task which can then be used to target potential customers for a marketing campaign.

The target variable consists of either 0 or 1, where 0 represents the purchase of a single vehicle by the customer, while 1 represents the purchase of multiple vehicles. As we aim to predict only these two possible classes outcome, this Binary Classification model is the right choice along with learning objective purpose.

Hyperparameters Selected.

Hyperparameters Name	Values Tested		
penalty:	elasticnet		
l1_ratio:	0.5	0 (l2)	1 (l1)
class_weight:	balanced		
solver:	saga		

\* Started with default hyperparameters for the Logistic Regression algorithm and assessed its performance and then adjustments to the hyperparameters in subsequent exercises to achieve optimal performance.

\* To address the issue of false negative errors, applied regularization techniques and evaluated the model's performance to check if it can help minimize the misclassification of potential customers. Specifying l1\_ratio=0, is same as utilizing penalty='l2', and l1\_ratio=1, is equivalent to utilizing penalty='l1'.



	<p>* Utilising class_weight='balanced' hyperparameter to adjust the weights of the classes based on their proportion in the training set, which can help to address the issue of high imbalance classes and improve model performance.</p> <p>* While the solver='saga' hyperparameter is used to specify the algorithm used to optimize the model (also useful for large datasets or datasets with a high number of features).</p> <p>* Therefore, using class_weight='balanced' and solver='saga' can help to address issues related to class imbalance and model optimization, respectively, which can ultimately lead to better performance of the model.</p> <p>Because of limitations on time, it was not feasible to test other hyperparameters such as fit_intercept, intercept_scaling, and multi_class='ovr' of the Logistic Regression Classifier algorithm. However, would like to examine and test these hyperparameters in future experiments.</p>
--	--

3. EXPERIMENT RESULTS																		
Analyse in detail the results achieved from this experiment from a technical and business perspective. Not only report performance metrics results but also any interpretation on model features, incorrect results, risks identified.																		
3.a. Technical Performance	Score of the relevant performance metric(s). Provide analysis on the main underperforming cases/observations and potential root causes.																	
	<u>Performance Metrics</u>																	
	To evaluate the model’s performance utilizing the below performance metrics.																	
	<div><div>-</div>Precision</div>																	
	<div><div>-</div>Recall</div>																	
	<div><div>-</div>Weighted F1 Score</div>																	
<div><div>-</div>Binary F1 Score</div>																		
<div><div>-</div>Confusion Matrix</div>																		
<u>Baseline Performance (Null Accuracy)</u>																		
weighted F1 Score: 0.9608																		
<b>Exercise 1:</b>																		
Algorithm: Logistic Regression Classifier model																		
Hyperparameters: Default Hyperparameters																		
Performance Metrics:																		
<table><tr><td>Dataset</td><td>Precision</td><td>Recall</td><td>weighted F1 Score</td><td>binary F1 Score</td></tr><tr><td>Training:</td><td>0.8327</td><td>0.2189</td><td>0.9721</td><td>0.3467</td></tr><tr><td>Validation:</td><td>0.8301</td><td>0.2340</td><td>0.9721</td><td>0.3651</td></tr></table>				Dataset	Precision	Recall	weighted F1 Score	binary F1 Score	Training:	0.8327	0.2189	0.9721	0.3467	Validation:	0.8301	0.2340	0.9721	0.3651
Dataset	Precision	Recall	weighted F1 Score	binary F1 Score														
Training:	0.8327	0.2189	0.9721	0.3467														
Validation:	0.8301	0.2340	0.9721	0.3651														

	Error Rate		Error Count	
Dataset	False Positive	False Negative	False Positive	False Negative
Training:	0.12%	2.05%	97	1723
Validation:	0.13%	2.06%	27	432

\* From the above, the weighted F1 score and Precision score of both the training and validation sets, the logistic regression classifier with default hyperparameters appears to be performing well overall. However, the Recall score of 0.2189 for training and 0.2340 for validation, indicates that the model may miss some potential customers who are likely to purchase a new vehicle.

\* Additionally, the False Negative error for both the sets training with 1723 and validation with 432, informs that a considerable number of customers who is likely to buy a new vehicle are being predicted otherwise by the model. This misclassification could result in missing out on these potential customers during the marketing campaign.

\* Therefore, to address the issue of false negative errors, applying regularization techniques and evaluating the model's performance to check if it can help minimize the misclassification of potential customers.

### Exercise 2:

Algorithm: Logistic Regression Classifier model

Hyperparameters: penalty='elasticnet', l1\_ratio=0.5, class\_weight='balanced', solver='saga'

Performance Metrics:

Dataset	Precision	Recall	weighted F1 Score	binary F1 Score
Training:	0.0920	0.8680	0.8494	0.1665
Validation:	0.0966	0.8812	0.8515	0.1741

	Error Rate		Error Count	
Dataset	False Positive	False Negative	False Positive	False Negative
Training:	22.46%	0.35%	18881	291
Validation:	22.11%	0.32%	4647	67

\* The logistic regression classifier with 'elasticnet' regularisation/penalty (and other hyperparameters class\_weight='balanced' and solver='saga') with Recall score of 0.8680 for training and 0.8812 for validation informs that the model is good at identifying customers who are likely to purchase a new car but may be missing out on many of them.

\* However, the binary F1 score has degraded for both the training with 0.1665 and validation with 0.1741, indicating that the model is not performing well in correctly classifying the positive class instances.

\* Additionally, the False Negative error has significantly dropped to 291 for training and 67 for the validation set, thereby increasing the Recall score but compromising the Precision score on both sets.

### Exercise 3:

Algorithm: Logistic Regression Classifier model

Hyperparameters: penalty='elasticnet', l1\_ratio=0, class\_weight='balanced', solver='saga'

Performance Metrics:

Dataset	Precision	Recall	weighted F1 Score	binary F1 Score
Training:	0.0920	0.8680	0.8494	0.1665
Validation:	0.0966	0.8812	0.8515	0.1741

	Error Rate		Error Count	
Dataset	False Positive	False Negative	False Positive	False Negative
Training:	22.46%	0.35%	18880	291
Validation:	22.11%	0.32%	4647	67

\* By utilizing regularization and implementing an 'l2' penalty along with other hyperparameters such as `class_weight='balanced'` and `solver='saga'`, the performance is identical to the previous model with elasticnet penalty.

\* Even the False Negative error remain same with 291 for training and 67 for the validation set.

#### Exercise 4:

Algorithm: Logistic Regression Classifier model

Hyperparameters: `penalty='elasticnet'`, `l1_ratio=1`, `class_weight='balanced'`, `solver='saga'`

Performance Metrics:

Dataset	Precision	Recall	weighted F1 Score	binary F1 Score
Training:	0.0920	0.8680	0.8494	0.1665
Validation:	0.0966	0.8812	0.8515	0.1741
Testing:	0.0989	0.8561	0.8486	0.1774

	Error Rate		Error Count	
Dataset	False Positive	False Negative	False Positive	False Negative
Training:	22.46%	0.35%	18881	291
Validation:	22.11%	0.32%	4647	67
Testing:	22.28%	0.41%	5853	108

\* The performance metrics for Precision, Recall, weighted and binary F1 Scores are the same across all three regularization methods l1, l2, and elasticnet.

\* Regularization has led to a significant reduction in False Negative errors with 291 for training and 67 for validation set, resulting in an increase in the Recall score for both datasets.

\* However, binary F1 Score has significantly dropped indicating that the model is not performing well in identifying customers who are likely to purchase a new vehicle.

### 3.b. Business Impact

Interpret the results of the experiments related to the business objective set earlier. Estimate the impacts of the incorrect results for the business (some results may have more impact compared to others)

The logistic regression algorithm with default hyperparameters (exercise 1) had a Recall score of 0.2189 for training and 0.2340 for validation, indicating that the model may miss some potential

	<p>customers who are likely to purchase a new vehicle along with a huge number of False Negative errors.</p> <p>Applying regularization techniques such as l1, l2 and elasticnet, has led to a significant reduction in False Negative errors representing some customers who are likely to purchase new vehicles are being overlooked, with 291 for training, 67 for validation and 108 for the testing set, as opposed to the initial model which had 1723 for training and 432 for validation sets.</p> <p>However, the regularization has significantly dropped binary F1 Score, indicating that the model is not performing well in identifying customers who are likely to purchase a new vehicle which could lead to the loss of chances to make sales and earn revenue, and may also result in failure to achieve marketing campaign goals or objectives.</p> <p>Therefore, the Logistic Regression Classifier models from these 4 exercises, are not sufficiently generalized to accurately predict customers who are likely to purchase a new vehicle.</p>
<b>3.c. Encountered Issues</b>	<p>List all the issues you faced during the experiments (solved and unsolved). Present solutions or workarounds for overcoming them. Highlight also the issues that may have to be dealt with in future experiments.</p> <p>It was challenging to choose the suitable hyperparameters for this binary classification problem from all the hyperparameters options available for the Logistic Regression Classifier. So, I went through sklearn's documentation to understand the hyperparameters and selected the ones used in the 4 exercises.</p> <p>Although each hyperparameter holds its own specific significance, it can have varying effects on the performance of the model. Therefore, it is crucial to apply hyperparameters with caution, taking into consideration the problem being addressed and the characteristics of the dataset.</p>

<b>4. FUTURE EXPERIMENT</b>	
<p>Reflect on the experiment and highlight the key information/insights you gained from it that are valuable for the overall project objectives from a technical and business perspective.</p>	
<b>4.a. Key Learning</b>	<p>Reflect on the outcome of the experiment and list the new insights you gained from it. Provide rationale for pursuing more experimentation with the current approach or call out if you think it is a dead end.</p> <p>Both the logistic regression models, one with default hyperparameters and the other with regularization techniques, are not sufficiently generalized and show relatively inadequate overall performance. So, it would be beneficial to explore other binary classification algorithms that can better identify potential customers who are likely to purchase a new vehicle.</p> <p>Additionally, before calling it the dead end and explore other algorithms, would like to experiment and test other hyperparameters such as fit_intercept, intercept_scaling, and multi_class='ovr' to check if the model's ability to generalize could be enhanced and performance could be improved.</p>

#### 4.b. Suggestions / Recommendations

Given the results achieved and the overall objective of the project, list the potential next steps and experiments. For each of them assess the expected uplift or gains and rank them accordingly. If the experiment achieved the required outcome for the business, recommend the steps to deploy this solution into production.

I would carry out these two tasks.

- \* Train the Logistic Regression classifier algorithm with additional hyperparameters like `fit_intercept`, `intercept_scaling`, and `multi_class='ovr'` and review if the performance is improving.
- \* Analyse the model's performance using cross-validation technique.

To successfully deploy a binary classification algorithm in a production environment, it is advisable to carry out below steps.

- Scale and transform the model to handle large datasets.
- Select an appropriate deployment environment whether cloud-based or on-premises.
- Modify the model to suit production settings while complying with various security, ethical, and privacy guidelines.
- Conduct testing and monitoring of the deployed model.
- Periodically updated the model with new data.
- Provided documentation for usage.
- Review the performance of the model and retrain the model as needed.