

# EXPERIMENT REPORT [ B ]

## EXPERIMENT BACKGROUND

Student Name	Nathan Collins
Project Name	MLAA Assignment 2
Date	Apr 28 by 23:59
Deliverables	<MLAA_Assignment_2_2> <Experiment_Report_2>

### a. Business Objective

The overarching aim of the project is to utilise analytical and statistical methods to determine the likelihood that an existing customer of an automotive manufacturer will purchase a new vehicle. To achieve this, several experiments will be conducted using a provided historical dataset. The success of the model will empower business stakeholders to implement a cost-effective re-purchase campaign targeting existing customers who are prospective in purchasing a new or secondary vehicle.

The accuracy of the results will determine the value of the leads attained following the marketing campaign. Precise results may yield a positive return on investment by identifying these fulfilling customer leads, while incorrect results could lead to revenue loss through an unsuccessful marketing campaign, overstocking incorrect car varieties (models and segments) in anticipation of specific buyers, or granting unnecessary warranty and servicing inclusions with purchases, enacting a toll on the business, and damaging its longevity.

### b. Hypothesis

#### Alternative hypothesis:

There is a relationship between some features in the car sales dataset and the likelihood that existing customers of an automotive manufacturer will purchase a new vehicle.

#### Null hypothesis:

There is no relationship between the features in the car sales dataset and the likelihood that existing customers of an automotive manufacturer will purchase a new vehicle.

The null hypothesis assumes no relationship is present between the dataset features, while the alternative hypothesis suggests that there is a relationship. By comparing the model's performance against the null hypothesis, it determines whether the model is significantly better than random chance at predicting customer purchase behaviour.

Investigating these questions will assist with accomplishing the business goal by determining which existing customers who have purchased a second vehicle share certain features.

### c. Experiment Objective

Experiment 2. will consist of a **multivariate logistic regression** model, examining the **Target** variable (which determines if a customer has purchased more than 1 vehicle) against all other variables independently, each found in two separate data frames:

*Table 1 Data frames applied for analysis and modelling.*

Data Frame	Feature
<b>cars_ALL</b>	"all features" 18289 entries, 42 features
<b>cars_NAG</b>	"no age-gender" 128611 entries, 34 features

Working with two data frames facilitates a broader scope of outcome possibilities while retaining as many variables as possible without resorting to oversampling. As features such as **age\_band** and **gender** offer decisive analytical data and insights into the classification of certain population groups, they will be retained where possible.

For this experiment to be deemed a success, it is anticipated that precision and recall values will result in high coefficients, close to a value of 1. If a low or no depicted coefficient value is produced, it will pave the way for further experimentation with more complex algorithms.

These features will undergo further classification through evaluating model output coefficients to determine **feature importance**, where key traits about existing customers can be utilised to determine if they have purchased a second vehicle, achieving the business objective.

EXPERIMENT DETAILS

a. Data Preparation

< Data preparation method unchanged from Experiment 1. >

Data Understanding

The data frame was explored to derive foundational insights, indicating the initial seventeen features, two of which contain NaN values (**age\_band** and **gender**) and four consisting of object variables (**age\_band**, **gender**, **car\_model** and **car\_segment**) which will need to be transformed. All remaining features are full and are integer values.

Data Cleaning

In order to prepare the data for analysis and modelling, the ID column was removed to reveal 2726 duplicate entries. As duplicate entries can weigh certain outcomes and impede data accuracy and consistency, these were dropped.

Next, the NaN values were totalled: 109668 in **age\_band** (85.2% of total rows) and 67455 in **gender** (52.4% of total rows) and visualised with **missingno**. As age and gender offer alternative means of classification over the dataset, aside from car-related data, it was decided to retain these features in a separate yet more condensed dataset.

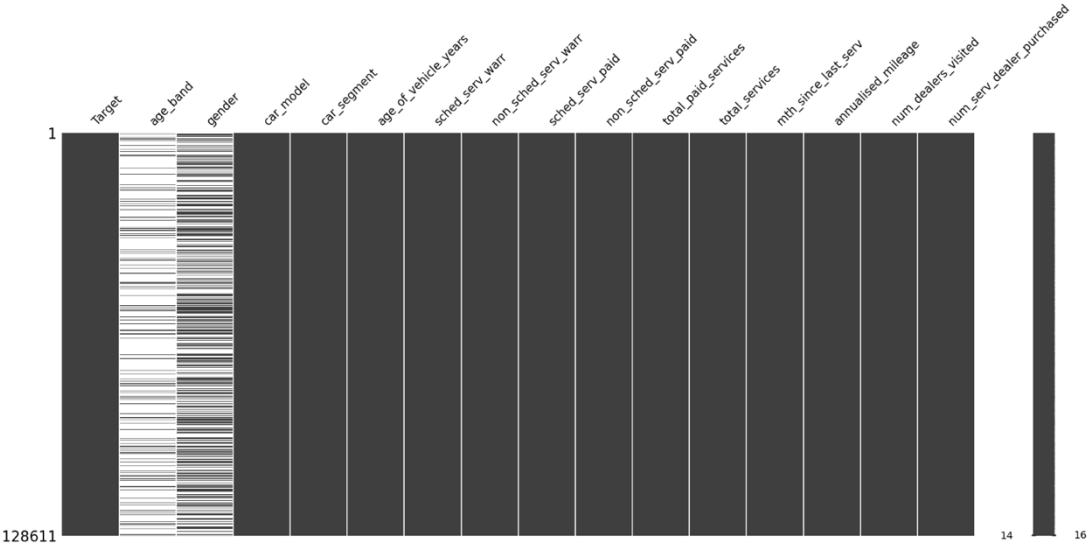


Figure 1 NaN values visualised with Missingno.

## Producing Two Data Frames with All Values Converted to Integers:

**cars\_ALL** and **cars\_NAG** (see Table 1.)

**For cars\_ALL:**

The **age\_band** distribution was visualised prior to conducting **one-hot-encoding**, converting the feature into six rows with integer values of 1 and 0.

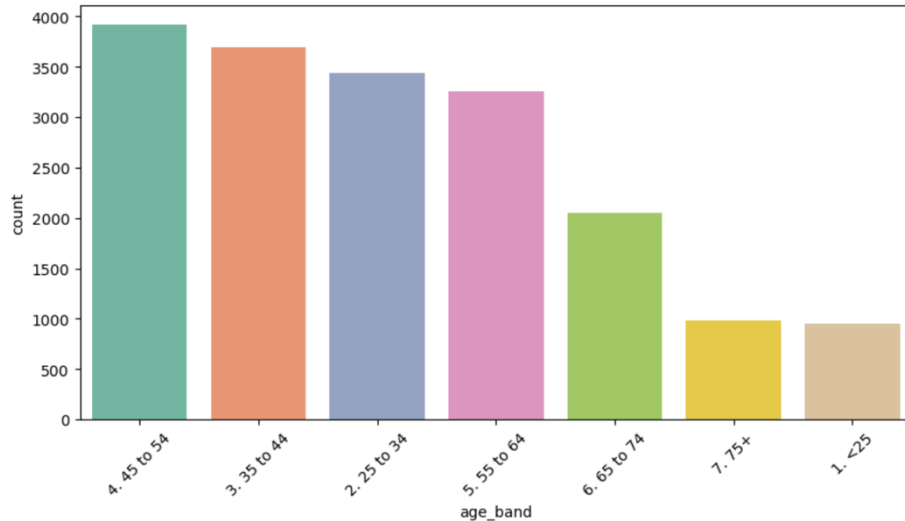


Figure 2 The **age\_band** feature visualised.

One-hot-encoding was subsequently performed on the **gender** feature, converting the variable into two rows with integer values of 1 and 0, male and female.

**For cars\_ALL and cars\_NAG:**

The **car\_segment** distribution was also visualised across both datasets prior to **one-hot-encoding**, converting the feature into four separate features with integer values of 1 and 0, Small/Medium, Large/SUV, LCV and Other

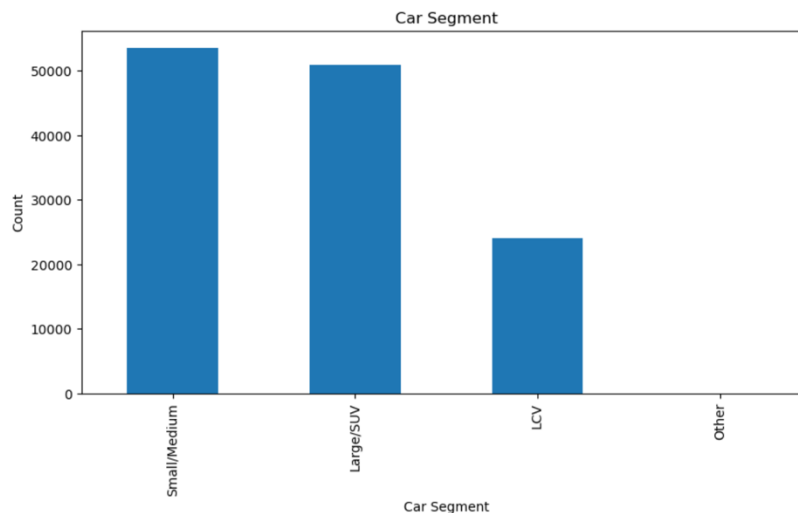


Figure 3 The **car\_segment** feature of the **cars\_NAG** data frame, visualised.

The **car\_model** feature was the last categorical variable converted to an integer feature through **one-hot-encoding**, yielding nineteen features representing different car models with integer values of 1 and 0.

The ranges of each dataset were visualised as a final measure to ensure no outliers or abstract inputs were present.

## Exploratory Data Analysis

Despite the problem residing in binary classification, linear correlations of all features against all other features of the dataset were first visualised using a heatmap, where the top 10 correlating features for each dataset were collated in a table.

	cars_ALL	correlation_ALL		cars_NAG	correlation_NAG
0	Target	1.000000		Target	1.000000
1	Male	0.033980		Large/SUV	0.015211
2	Large/SUV	0.027986		LCV	0.010342
3	4. 45 to 54	0.016244		car_model	0.000575
4	5. 55 to 64	0.013153		Other	-0.001319
5	car_model	0.008100		Small/Medium	-0.023228
6	LCV	0.004723		non_sched_serv_paid	-0.033297
7	7. 75+	-0.001011		num_dealers_visited	-0.053589
8	Other	-0.001083		num_serv_dealer_purchased	-0.058963
9	3. 35 to 44	-0.004205		annualised_mileage	-0.080251
10	2. 25 to 34	-0.007470		non_sched_serv_warr	-0.088442

Figure 4 The cars\_ALL dataset and the cars\_NAG dataset highest correlating features.

Next, the percentage of the **Target** variable class of interest (integer values of 1) was visualised with a pie graph. Indicating that only 2.7% of the 128611 entries were the population group of interest. This interpretation establishes the dataset as unbalanced.

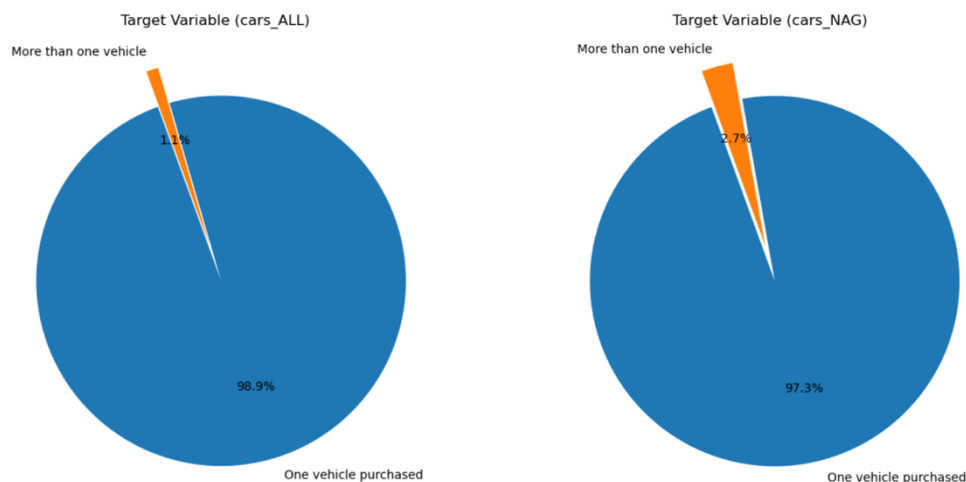


Figure 5 Pie graph visualisations of the Target variable and the class of interest.

Further visuals were subsequently constructed to gauge the perceived influence of the **gender** and **age\_band** variables against other key classification features, such as the total, which share the **Target** variable and the relative **car\_segment** and **car\_model** choices made with their second car purchases. Overall it was interpreted that males tended to purchase a second car more than females, with the largest age brackets resting between 45 to 54. Out of these groups, the highest interest resided with Large/SUVs, models 5 and 3, among males. It is important to note that these features represent marginal percentages of the total and are not a full representation of the target cohort.

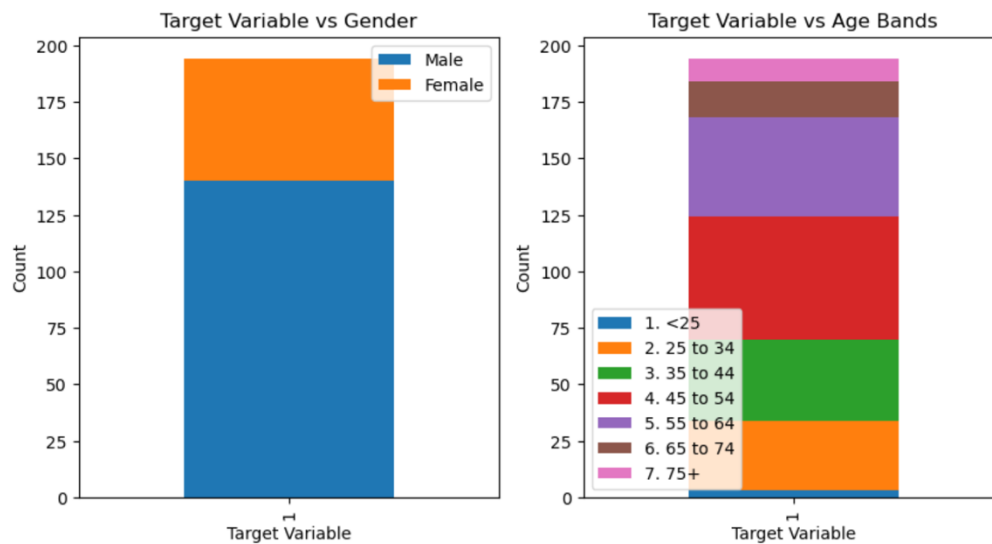


Figure 6 A segmented bar graph illustrating the male and female cohorts, beside the age band cohorts, with a Target value of 1.

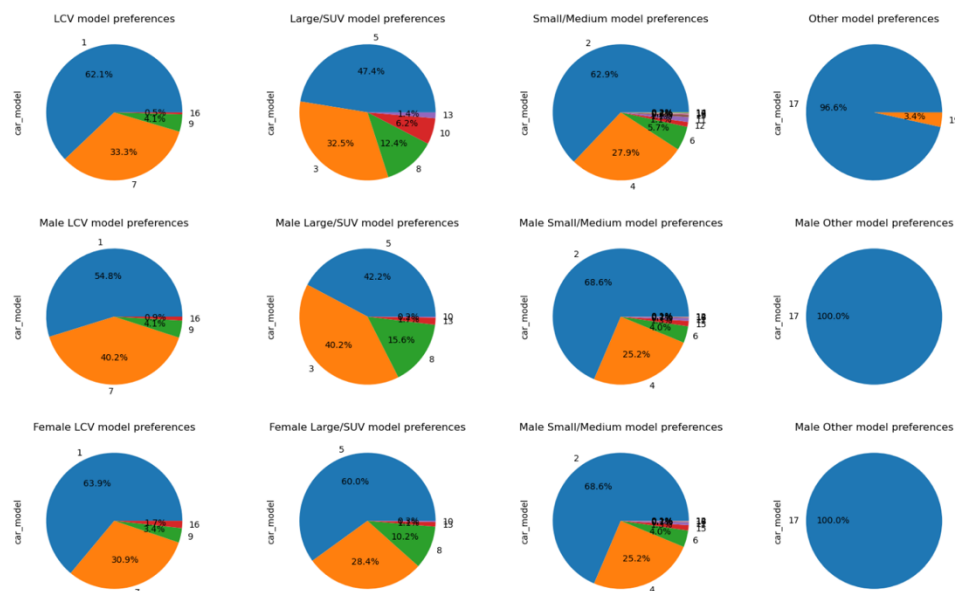


Figure 7 A series of pie charts illustrating the popularity of car models against the variety of car segments. The first row covers the overall majority, while the second and third rows cover the male and female preferences.

The final key exploratory analysis involved charting the frequencies of the **Target** variable with a value of 1 against **total\_services**, **annualised\_mileage**, **age\_of\_vehicle\_years**, and **non\_sched\_serv\_war**. The value was to visualise whether specific features of a customer's existing had an influence on their decision to purchase a new one. A key insight derived from this is that recurring customers typically engage with 2 to 3 total services (deciles), with slightly more non-scheduled services.

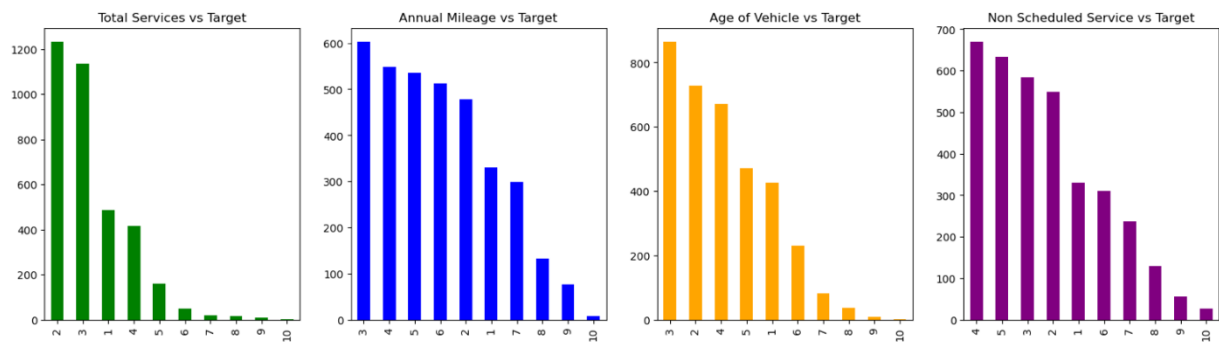


Figure 8 Four bar graphs illustrating *total\_services*, *annualised\_mileage*, *age\_of\_vehicle\_years*, and *non\_sched\_serv\_war* in deciles, against the target variable with a value of 1.

## b. Feature Engineering

< Aside from applying one-hot-encoding on **age\_band**, **gender**, **car\_model** and **car\_segment**, no further feature engineering was conducted for Experiment 2. >

## c. Modelling

### Selecting a performance metric

As the business objective resides with a binary classification problem, **Precision** was selected as the key performance metric to determine the model's success. Precision examines the proportion of true positives among the total positive predictions. By applying precision, a reduction in predicting the wrong customer is achieved. Should precision metrics share output similarity across models, an additional metric, **Recall**, will also be evaluated. Recall examines the proportion of true positives among all actual positive outcomes, meaning higher recall results in more lost opportunities for the marketing campaign.

Moving forward, a **precision-recall curve** will be constructed with each model to act as a visual depiction of the model's performance. In a successful model, the concavity of the curve determines the overall accuracy of results. In addition, an F1 and accuracy score will be shown alongside a confusion matrix. The variation in reporting performance metrics facilitates a deeper understanding of the model's predictive capability and how it compares with future experiments

### Establishing the model:

All models across all experiments, will be constructed through a Python function. This helps establish consistency across all metrics within the dataset and throughout the transformation of the data prior to and during machine learning. The subsequent model selected for this stage of the experiment is intended to expand upon Experiment 1's insights and tangent away from a simple model.

### Multivariate Logistic Regression

**Complexity:** Intermediate

$$\ln\left(\frac{\hat{p}}{(1-\hat{p})}\right) = b_0 + b_1X_1 + b_2X_2 + \dots + b_pX_p$$

Prior to constructing the function, both data frames (cars\_ALL and cars\_NAG) were first split at a **random\_state** of 100 into a training and testing set with embedded stratification (operating under the conventional **80% training / 20% testing** approach).

A function for the multivariate logistic regression was then defined as **ML\_results**, where subsequent iterations of the function could interchange the **dataset**, **feature**, **target**, and **model** variables. The function would first normalise the data first through **StandardScaler()**, which subtracts the mean from each data point and divides it by the standard deviation. The following enables data points with different units of measurement to be compared against one another and is essential prior to running a multivariate model. The output variables from this process produces **X\_train\_s** and **X\_test\_s**.

The function then instantiates and fits the model to the training data (**.fit**) and creates predictions on the training and test data (**.predict**).



To gauge the model's performance, the function will print a confusion matrix, followed by the accuracy score of both training and testing sets, followed by an F1 score of both training and testing sets (**confusion\_matrix**, **accuracy\_score**, **f1\_score**).

The last aspect of the function is establishing a Boolean with **.predict\_proba** to plot a precision-recall curve (**.precision\_recall\_curve**). The precision results extracted to construct the curve will also be embedded in a Pandas data frame displayed below, presenting all precision scores of **0.75** or **greater**.

**No baseline** metric (for assessing null accuracy) was applied to compare performance against naïve predictions and determine whether the model is adding value. This was selected, as the majority of the data in the **Target** cohort is negative, making the null accuracy equal to the proportion of negative samples.

### Hyperparameter tuning:

Hyperparameter tuning comprises determining the most fitting set of hyperparameters for a model. These parameters are not learned from the data and are required to be set by the user prior to training a model. For this experiment, a grid approach (**.GridSearchCV**) will be applied to a hyperparameter dictionary (**hyperparameters\_dict**). The dictionary will be embedded into a different function with the same functionality as **ML\_results**, although renamed to **ML\_results\_cv**. Should hyperparameters need to be tuned, the function **ML\_results\_cv** will be applied in conjunction with separate hyperparameter dictionaries.

These finalised functions were first applied to the **cars\_ALL** and **cars\_NAG** data frames where appropriate, defining the **dataset**, **feature**, **Target** and **model** with each iteration.

### Models to consider for future experiments:

#### Support Vector Classification

- Kernel functionality for higher dimensional class separation.
- Strong when boundaries between classes are clear.
- Capable of handling large datasets.

#### Random Forest Classification

- Reduced overfitting functionality through decision trees.
- Accounts for missing values, outliers, and nonlinearity.
- Wide range of hyperparameter tweaking possibilities.

#### XGBoosting Classification

- Compounds simpler models to build a larger analysis.
- Capable of accounting for missing data.
- Requires additional data preparation and parameter tuning than previous algorithms.

## EXPERIMENT RESULTS

### a. Technical Performance

#### Evaluation of precision

In general, a high precision score indicates the model has a low rate of false positives, while a high recall score indicates that the model has a low rate of false negatives. To meet the business objective, successful model outcomes should aim to achieve both precision and recall scores close to 1.

#### Technical evaluation

The first multivariate logistic regression was performed on the cars\_ALL dataset without hyperparameter tuning, resulting in a moderately high precision but low recall. The **precision** score of 0.7 means that out of all the instances predicted as positive, 70% were actually positive, while the remaining 30% were false positives.

The **recall** score of 0.3, however, means that out of all the instances that were actually positive, only 30% were correctly identified by the model, while the remaining 70% were false negatives. With such a large low recall rate, the model is presently not ideal for classifying the cars\_ALL dataset.

The precision-recall curve shape additionally lacks concavity and uniformity, indicating “noise” is apparent within current features and the possible need to tune hyperparameters.

From a business perspective, the model in its current form is not adequate enough to accurately meet the business objective in forecasting the features that relate to the target class of prospective buyers.

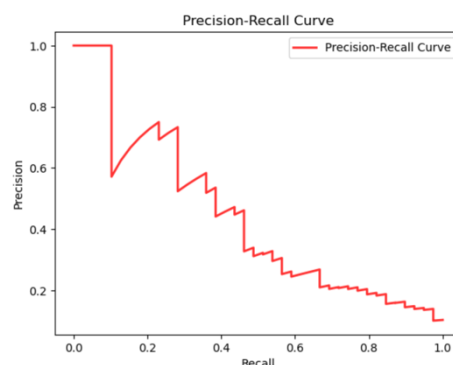
#### Performance Metrics for Multivariate Logistic Regression 1.

**Dataset:** cars\_ALL

Confusion Matrix Training Set	
14457	19
108	47
Confusion Matrix Testing Set	
3615	3615
30	30

Training Set	
Precision	0.712
Recall	0.303
Testing Set	
Precision	0.692
Recall	0.231

<b>Accuracy</b> Training Set	0.99
<b>Accuracy</b> Testing Set	0.99
<b>F1 Score</b> Training Set	0.99
<b>F1 Score</b> Testing Set	0.99



The second multivariate logistic regression was performed on the cars\_NAG dataset without hyperparameter tuning, resulting in, likewise, a moderately high precision but low recall. The **precision** score of 0.8 represents that out of all the instances predicted as positive, 80% were actually positive, while the remaining 20% were false positives.

The **recall** score of 0.2, however, means that out of all the instances that were actually positive, only 20% were correctly identified by the model, while the remaining 80% were false negatives. With such a low recall rate, the model is presently not ideal for classifying the cars\_NAG dataset.

The precision-recall curve, while more uniform than the previous curve, lacks concavity, indicating a possible lack of relationship between features and the possible need to tune hyperparameters.

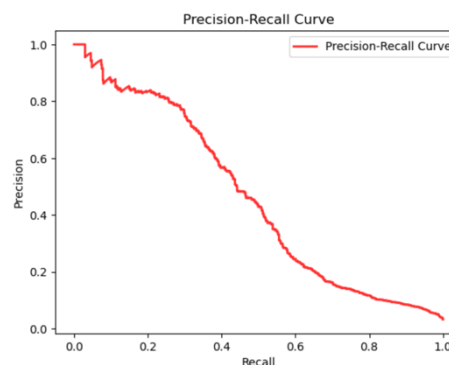
From a business perspective, the model in its current form is not adequate enough to accurately meet the business objective in forecasting the features that may relate to the target class of prospective buyers.

## Performance Metrics for Multivariate Logistic Regression 2. Dataset: cars\_NAG

Confusion Matrix Training Set	
99951	120
2193	624
Confusion Matrix Testing Set	
24991	28
563	141

Training Set	
Precision	0.839
Recall	0.222
Testing Set	
Precision	0.834
Recall	0.200

<b>Accuracy</b> Training Set	0.98
<b>Accuracy</b> Testing Set	0.98
<b>F1 Score</b> Training Set	0.97
<b>F1 Score</b> Testing Set	0.97



The third multivariate logistic regression was also performed on the cars\_NAG dataset with hyperparameter tuning. These parameters were tweaked across several iterations, altering the combinations of **penalty**, **C** and **solver**. The most successful of these resulted in virtually unchanged performance metrics from the previous model without hyperparameter tuning. The **precision** score of 0.8 represents that out of all the instances predicted as positive, 80% were actually positive, while the remaining 20% were false positives.

The **recall** score of 0.2, however, means that out of all the instances that were actually positive, only 20% were correctly identified by the model, while the remaining 80% were false negatives. With such a large low recall rate, the model is still not ideal for classifying the cars\_NAG dataset.

The precision-recall curve, while more uniform than the previous, still lacks concavity, indicating a possible lack of relationship between features with the current model applied.

From a business perspective, the model in its current form is not adequate enough to accurately meet the business objective in forecasting the features that may relate to the target class of prospective buyers.

### Performance Metrics for Multivariate Logistic Regression 3.

**Dataset:** cars\_NAG

**Hyperparameters tuned:** 'penalty' (l1, l2, elasticnet, none)

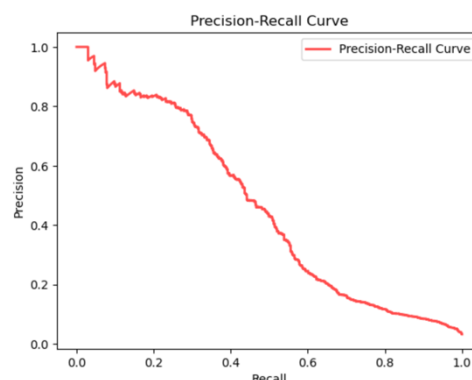
'C' (0.1, 0.5, 1.0)

(all similar outcomes) 'solver' (newton-cg, lbfgs, liblinear, sag, saga)

Confusion Matrix Training Set	
99951	120
2193	624
Confusion Matrix Testing Set	
24991	28
563	141

Training Set	
Precision	0.839
Recall	0.222
Testing Set	
Precision	0.834
Recall	0.200

<b>Accuracy</b> Training Set	0.98
<b>Accuracy</b> Testing Set	0.98
<b>F1 Score</b> Training Set	0.97
<b>F1 Score</b> Testing Set	0.97



## Feature Significance

By ranking the coefficient values generated from the model, the significance of each feature was charted. As most features displaying higher coefficient values are car\_model varieties and stem from poor precision-recall curves, it is unlikely these features provide an impactful prediction utility for aligning with the business objective.

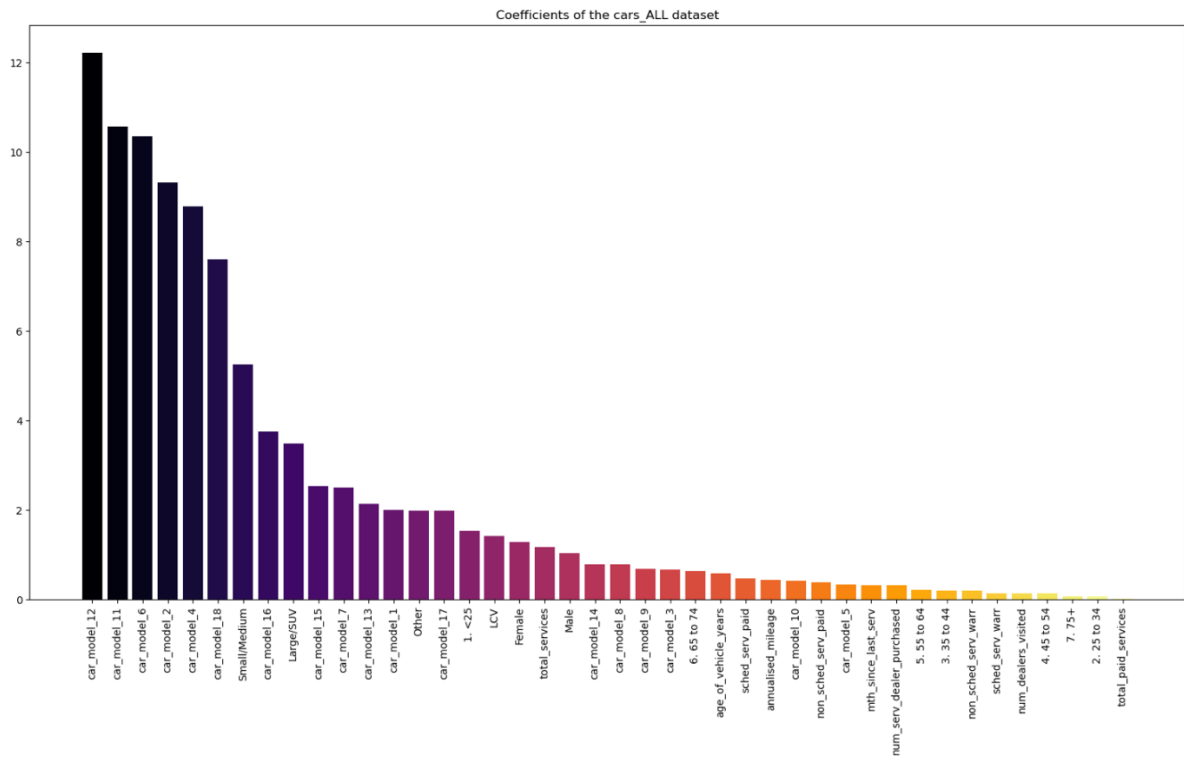


Figure 9 Coefficient values of the cars\_ALL dataset.

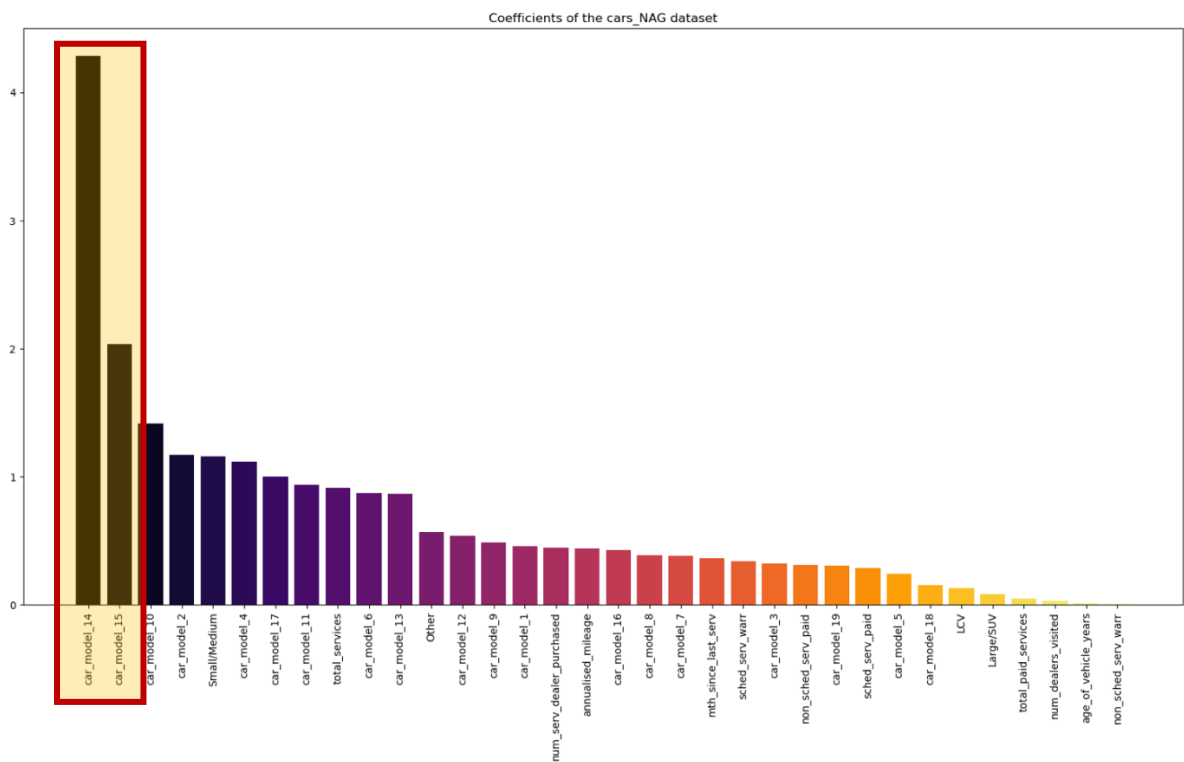


Figure 10 Coefficient values of the cars\_NAG dataset.

## b. Business Impact

As the results of each **multivariate logistic regression** model bear no significant predictability capability about the relationship between the **Target** variable and any other variable in the corresponding dataset, it is unfortunately deemed an unsuccessful experiment and provides no input towards meeting the business objective. This data will ultimately not aid in evaluating existing customers to procure leads for marketing.

The results are verified to be accurate by examining the **confusion matrix**, **F1\_score** and **Accuracy\_score**. Despite this, as this experiment model will not be provided to the business, there is no concern over providing inaccurate data and the ramifications which would follow.

## c. Encountered Issues

### The application of both data frames

Following the initial exploration of the missing values present within the **age\_band** and **gender** features, deliberation and evaluation of their value ensued, as the value in retaining these classification features needed to be considered. While the decision to preserve these features in a separate data frame was enacted, it ultimately scaled the workload. The use of the cars\_ALL data frame was ultimately discontinued during this experiment throughout hyperparameter tuning, as the resulting precision-recall curve was jagged and lacked uniformity.

### Coding learning curves and mitigating human error

After concluding that the most time-effective method for meeting the business goal was through constructing a classification function with an embedded splitting, model fitting and scoring capacity, establishing an approach occupied days of my time. It is anticipated that with practice, this issue will be avoided for future projects. By consulting library descriptions, forums and advisory web pages, the solution was found by interchanging the **dataset**, **feature**, **target**, and **model** with each function iteration. This issue is fortunately lessened with future experiments, as understanding the foundation to build future models was a priority.

Overcoming human error likewise presented itself as an inconvenience. As the functions are quite large, duplicating each and interchanging the variables inputted resulted in some initial incorrect results, remedied through constant revision and repetition.

### Resolving time limitations

While the processing time for these models was relatively efficient, the trial-and-error process in tweaking hyperparameters proved difficult to balance. This was resolved by integrating the **search.best\_estimator\_** feature, which enabled predictions to be made on better hyperparameter choices with each subsequent iteration.

## FUTURE EXPERIMENT

### a. Key Learning

#### Exploratory Data Analysis

*< Data Insights from EDA remain unchanged from Experiment 1. >*

#### Multivariate Logistic Regression

- Overall, progress has been made in the right direction. As it stands, no key insights into business decisions can be enacted, as performance metrics (despite a moderately high precision, still yielded a low recallability). The following indicates the need to explore other models.
- The implementation and application of hyperparameter tuning is a time-consuming process that resulted in no discernible differences in this model with these datasets.
- The significant features evaluated from this model are mostly car\_model varieties, which is an unlikely conclusion.

### b. Suggestions / Recommendations

Future experiments should explore further feature interactions in conjunction with further **multivariate classification** models. The project would benefit from exploring more complicated methods of manipulating and transforming data for meaningful interpretation. Future experiments may also benefit from further **feature engineering**, offering a means to transform data into more interpretable relationship modelling. As feature engineering can also reduce the dimensionality of the data and remove redundant aspects, it may lead to more promising precision and recall performance metric outcomes.

Other models to consider next include:

*< Directional insights remain unchanged from Experiment 1. >*

#### Support Vector Classification (SVC)

This classification algorithm may prove useful in future experiments, as it finds a hyperplane to maximally separate two classes, offering kernel functions like linear, polynomial, or radial options to better separate these classes.

Pros:

- Provides an effective means of separating non-linear datasets.
- Handles “high-dimensional” data.
- Regularisation functionality to prevent overfitting.

Cons:

- Slow for large datasets and computationally taxing.
- Sensitive to choice of kernel function and hyperparameters.

## Random Forest Classification

This classification algorithm may also provide value in future experiments by predicting classes of a new observation based on a highly adjustable set of input features.

### Pros:

- Offers functionality to handle outliers (not applicable) and missing data.
- Lower susceptibility to overfitting through bagging.
- Provides modelling functionality to a non-linear set of data.

### Cons:

- Requires a large set of trees to achieve optimal performance, which results in computationally intensive processing requirements.
- As the existing dataset is imbalanced, there is a poorer performance probability.

## XGBoost Classification

The XGB classification algorithm is a popular algorithm that combines multiple weak prediction models, like decision trees, iteratively adding trees and correcting errors of previous trees to generate an accurate prediction.

### Pros:

- Efficient processing times.
- Capable of handling **imbalanced** and missing data.
- Scalable for larger datasets.

### Cons:

- Typically requires more data preparation and parameter tuning than other algorithms.
  - Typically doesn't perform well with too many categorical features.
  - Prone to overfitting.
- 
- XGBoost is an optimised gradient-boosting algorithm that builds an ensemble of weak models to make accurate predictions.
  - It can handle missing data and has excellent predictive power for high-dimensional data.
  - However, XGBoost may require more data preparation and parameter tuning than other algorithms and can be more prone to overfitting.

---

## Deployment

As the experiment did not achieve the required outcome for the business, deployment into production is not recommended at this current time.