

# EXPERIMENT REPORT

Student Name	Shaqrان Bin Saleh
Project Name	Classification Models
Date	26/4/2024
Deliverables	Notebook Name : 36106_AT2_25010238_experiment_5.ipynb model name : Extra Trees Classifier

## 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

The goal of this project is to build a predictive model with accurate likeliness of repurchasing a car for customers.

The model's predictions can be used in a myriad of ways. The model can be used to segment customers by repurchase likelihood, focusing marketing efforts on high-potential features of the dataset for efficient resource use. Insights from the data and trained model can also enhance customer retention strategies and optimize inventory management. Additionally, the results would provide feedback to refine sales tactics and customer interactions, continually improving approaches to meet customer needs effectively.

Accurate predictions would mean enhance efficiency in targeting the right customers in reducing marketing efforts and higher sales. Through proper data exploration and trained model insights proper understanding of customer needs and help improve customer satisfaction.

Inaccurate repurchase likelihood prediction by the model might result in diminished business opportunities, reduced ROI, wasted marketing budget, and unhappy clients. Customers may become disconnected from a brand due to improper targeting, which could be perceived as spam and harm the company's reputation.

### 1.b. Hypothesis

I wish to test the hypothesis that several characteristics:  
Gender, The model of vehicle, The type of vehicle, Age of their last vehicle, Number of scheduled services used under warranty, Number of non-scheduled services used under warranty, Amount paid for scheduled services, Amount paid for nonscheduled services, Amount paid in total for services, Total number of services, The number of months since the last service, Annualized vehicle mileage, Number of different dealers visited for servicing, Number of services had at the same dealer where the vehicle was purchased have significant effects on the customers' likeliness of repurchasing a vehicle.

I like to get the answers of few questions like "What are the primary repurchasing variables for customers and how do they interact to influence their repurchasing likelihood?" By testing this hypothesis and answering the accompanying question, we can acquire a deeper understanding of the factors that influence customers' likelihood of repurchasing.

1.c. Experiment Objective	<p>From this experiment we will be deploying a predictive model using Extra Trees Classifier model that accurately handles the data and make accurate predictions.</p> <p>The goal is to identifying key variable effecting the customers' likelihood of repurchasing, developing an accurate model, apply it on the dataset and find the accuracy of the model basing upon the different hyperparameter tuning.</p>
---------------------------	--

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	<p>For proper data preparation, the following steps were taken:</p> <ol style="list-style-type: none"> <li>1. The assigned .csv file was loaded into the notebook as a dataframe for the experiment.</li> <li>2. The dimensions (rows, columns) of the data loaded from the .csv file was found. This step was performed so that we can properly understand the amount of data we are working with. It also gives us information about the number of features present in the data.</li> <li>3. We have 17 features among which 1 is our target variable which is "Target". We explored all the features of the dataset.</li> <li>4. We check all the features if there are any null values and handle the situation.</li> <li>5. We check the statistics (count, mean, standard deviation, minimum values, maximum value) of the features.</li> <li>6. We drop the rows containing missing values from the dataset.</li> <li>7. We drop the rows containing duplicate data.</li> <li>8. Since we have 4 features containing categorical data. We introduced OneHotEncoder through which we transform the data and concatenate the new columns into our dataset. And finally remove the previous columns on which the transformation was done.</li> <li>9. The feature "Target" is our target variable so we separate it as our 'y' variable and the rest were are considering our independent variable which</li> <li>10. Then we split the data to three categories. Which are training data, testing data and validation data.</li> <li>11. We distribute the data firstly into two segments Train + Validation and Test sets. The training + validation set would hold 80% of the data and rest 20% is testing data.</li> <li>12. Then we again split the Train + Validation to separate Training and Validation sets. The validation would hold 30% of the Train + Validation data, meaning it hold 24% of the original data. The rest of the 56% of the original data is reserved for the training dataset.</li> <li>13. And finally, we performed feature scaling</li> </ol> <p>All the above steps were vital in this experiment is also vital for the next experiments.</p>

<p>2.b. Feature Engineering</p>	<p>While performing data preparation, we found out that the features : 'age_band', 'gender', 'car_model', 'car_segment' are categorical and are all important in the influence of the target variable. So, we had to encode these features. We achieved this encoding using the OneHotEncoder()</p> <p>As we concatenated the new columns, we made sure to drop the columns on which the transformation was performed.</p>
<p>2.c. Modelling</p>	<p>For the experiment to predict customer repurchase likelihood, the Extra Trees Classifier model was chosen.</p> <p>Extra Trees Classifier is straightforward to implement and understand.</p> <ul style="list-style-type: none"> <li>• Extra Trees builds a "forest" of numerous decision trees. Crucially, these trees are constructed with randomness, meaning each tree uses different features and random split points. This diversity prevents overreliance on specific patterns in the training data and allows the model to adapt better to unseen data</li> <li>• Traditional decision trees can overfit the training data, memorizing noise or irrelevant details. Extra Trees combats this by injecting randomness into the tree-building process</li> <li>• Extra Trees leverages the "random forest" concept to provide valuable insights into feature importance. By analyzing how each feature contributes to the final prediction across the entire ensemble of trees, it identifies which features are most influential in predicting the target variable.</li> </ul> <p>The decision to use Extra Trees Classifier should be aligned with the specific characteristics of the data and the business or research objectives.</p> <p>While performing the experiment hyperparameters tuning was done. We took the help of Random Search CV. Through its estimator we found out the best values for the hyperparameters.</p> <p>The hyperparameter which were tuned were: n_estimator, max_depth, min_samples_leaf, max_features. The best values for the hyperparamters was:</p> <p>The n_estimator: 155  The max_depth: 6  The min_samples_leaf: 2  The max_features: 2</p> <p>The hyperparameters we tuned for your Extra Trees model are crucial for controlling the model's complexity and its ability to generalize to new data</p> <ul style="list-style-type: none"> <li>• The hyperparameter n_estimator specifies the number of decision trees included in the random forest model. While increasing the number of trees generally enhances the model's stability and accuracy, it comes at the expense of greater processing time.</li> <li>• The hyperparameter max_depth is basically the maximum depth of the tree. With limiting the depth of the tree we can prevent overfitting but again setting it too low will hamper the models ability to learn complex patterns.</li> <li>• The hyperparameter min_samples_leaf is the minimum number of samples a leaf node can have. There is a trade-off between the minimum number of samples required for a leaf node in a decision tree model and the model's generalization performance. A higher minimum number of samples reduces variance and can lead to improved ability to</li> </ul>

perform well on unseen data but may introduce bias. Conversely, a lower minimum allows the model to capture more detail from the training data but increases the risk of learning patterns specific to the training data that do not generalize well which is overfitting.

- The max\_features determines the number of features to consider when looking for the best split during the construction of each tree in the forest.
-

## 14. 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

The performance metrics used in this experiment are :

1. Accuracy Score
2. Precision Score
2. F1 Score
3. Confusion Matrix

We used the Extra Forest Classifier model and trained it multiple times as we performed multiple Hyperparameter tuning.

The results of the metrics according to different hyperparameters are as follow:

1. For default hyperparameters:

- Accuracy Scores:
  - a. On Training set: 0.9897560975609756
  - b. On Validation set: 0.9895311788802913
  - c. On Testing set: 0.9882577826324412
- F1 Scores:
  - a. On Training set: 0.0
  - b. On Validation set: 0.0
  - c. On Testing set: 0.0
- Precision Score:
  - a. On Training set: 0.0
  - b. On Validation set: 0.0
  - c. On Testing set: 0.0
- Confusion Matrix:  
[[3619 0]  
[ 43 0]]

2. For the best hyperparameters found from Grid Search ('max\_depth': 6, 'max\_features': 2, 'min\_samples\_leaf': 2, 'n\_estimators': 155 ):

- Accuracy Scores:
  - a. On Training set: 0.9897560975609756
  - b. On Validation set: 0.9895311788802913
  - c. On Testing set: 0.9882577826324412
- F1 Scores:
  - a. On Training set: 0.0
  - b. On Validation set: 0.0
  - c. On Testing set: 0.0
- Precision Score:
  - a. On Training set: 0.0
  - b. On Validation set: 0.0
  - c. On Testing set: 0.0
- Confusion Matrix:  
[[3619 0]  
[ 43 0]]

	<p>Main Issues and Causes:</p> <ul style="list-style-type: none"> <li>• Our analysis revealed a concerning discrepancy between the model's overall accuracy and its performance on the minority class. While the model achieved high accuracy across all data splits (training, validation, testing), the F1 score, precision, and recall scores for the minority class were all zero. This suggests an "accuracy paradox": the model excels at predicting the majority class but fails entirely at identifying instances of the minority class.</li> <li>• Further investigation through the confusion matrix confirmed this issue. Both the default and best hyperparameter settings resulted in a matrix with the following structure: [3619, 0; 43, 0]. This indicates a complete absence of true positives (correct predictions) for the minority class(actual repurchasing customers). All 43 minority class instances were misclassified as the majority class.</li> <li>• A significant class imbalance is evident within the data. This disparity can lead to a model biased towards the majority class. Such a model might simply predict the majority class for all inputs, neglecting the presence of the minority class altogether.</li> <li>• Despite the high accuracy reported, the model's complete inability to identify true positives raises concerns about overfitting. This suggests the model might not be capturing the underlying patterns crucial for class distinction. Overfitting on the majority class data is a potential culprit for this behavior.</li> </ul>
<p>3.b. Business Impact</p>	<p>According to our objective which we set earlier which was to deploy a predictive Extra Trees Classifier model that performs well. We can come the following conclusions about the results of this experiment:</p> <ul style="list-style-type: none"> <li>• The model exhibits a concerning bias towards the majority class. This is evident from the high overall accuracy coexisting with zero values for precision, recall, and F1-score on the minority class. While excelling at predicting frequent outcomes, the model fails to identify the less frequent, potentially more critical ones.</li> <li>• The employed Grid Search hyperparameter tuning did not address the class imbalance issue. Despite applying both default and tuned hyperparameters, the confusion matrices and metric scores for the minority class remained unchanged. This suggests that the tuning process did not tackle the underlying problem.</li> <li>• The consistent failure to detect the minority class compels to explore alternative training strategies. These might include: Assigning higher weights to the minority class instances during training can incentivize the model to focus on learning from these rarer examples. Techniques like oversampling or undersampling the data can help balance the class distribution and improve the model's performance on the minority class.</li> </ul> <p>And the impacts the business might face due to imperfect results would be:</p> <ul style="list-style-type: none"> <li>• If the model is deployed in this critical decision-making business, consistent inaccuracies or a failure to perform as expected could lead to loss of trust among stakeholders. Because even though perfect scores across all datasets might be initially impressive but are suspicious and could suggest overfitting or data leakage.</li> <li>• Missed opportunities due to incorrect model predictions can be significant. For instance, in a marketing campaign of this business, failing to correctly identify potential repurchasing customers (false negatives) could mean missed revenue opportunities.</li> </ul>

3.c. Encountered Issues	<p>As we performed this experiment, we were encountered with the following issues:</p> <ul style="list-style-type: none"> <li>• The data exhibits a significant class imbalance, with a considerably larger majority class compared to the minority class. This disparity led the model to prioritize predicting the majority class, resulting in zero precision, recall, and F1-score for the minority class.</li> <li>• The hyperparameter tuning strategy failed to address the issue of class imbalance. Both the default and best tuned hyperparameter settings provided identical confusion matrices and metric scores for the minority class. This implies that the tuning process was unable to rectify the model's core limitations in handling imbalanced data.</li> <li>• The primary reliance on accuracy as the evaluation metric misled the assessment of the model's performance. While the model achieved high accuracy by overwhelmingly predicting the majority class, it completely failed to identify the minority class correctly.</li> </ul>
-------------------------	--

1. FUTURE EXPERIMENT	
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.	
4.a. Key Learning	<p>Reflecting on the outcomes and performance of the Random Forest Classifier model experiments for predicting customer repurchase likelihood, several key insights emerge.</p> <ul style="list-style-type: none"> <li>• The experiment emphasized the significant influence of class imbalance on model behavior. When not addressed, even robust models like Extra Trees can become biased towards the majority class, neglecting the crucial minority class entirely.</li> <li>• While hyperparameter tuning is a valuable tool, it may have limitations in overcoming fundamental issues like data imbalance or inherent model biases.</li> <li>• This work underlined the potential pitfalls of relying solely on accuracy, especially for imbalanced datasets. The experiment demonstrated the need for more informative metrics like precision, recall, F1-score, and the confusion matrix to provide a comprehensive and realistic assessment of model performance across all classes.</li> </ul> <p>Given the insights and the issues encountered, here are the reasons to continue experimenting with the Extra Trees Classifier:</p> <ul style="list-style-type: none"> <li>• Compared to other tree-based models, Extra Trees exhibits a natural resistance to overfitting.</li> <li>• The model offers a multitude of hyperparameters (e.g., max_depth, min_samples_leaf, n_estimators) that remain open to further experimentation.</li> <li>• Extra Trees boasts computational efficiency, particularly when dealing with large datasets. This stems from its method of constructing trees using random splits, making it an attractive choice in such scenarios compared to many other ensemble methods.</li> </ul>

#### 4.b. Suggestions / Recommendations

The experimentation with the Extra Trees Classifier for predicting customer repurchase likelihood identified several pathways for further improvement. These refinements can be categorized into two main phases: model optimization and deployment preparation.

##### 1 Model Optimization:

- i) Moving beyond Grid Search, more sophisticated techniques like randomized search or Bayesian optimization can be employed to identify optimal settings for hyperparameters like max\_depth, min\_samples\_leaf, max\_features, and n\_estimators. These advanced algorithms explore the hyperparameter space more efficiently, potentially leading to superior model performance
- ii) Given the class imbalance issue, various strategies to balance the class distribution and improve minority class prediction should be investigated. These techniques include:
  - (1) Synthetic Minority Over-Sampling Technique (SMOTE) artificially generates synthetic data points for the minority class, effectively increasing its representation within the training data.
  - (2) Directly assigning higher weights to minority class instances during training incentivizes the model to prioritize learning from these rarer examples.
  - (3) Reducing the number of majority class instances can help counteract its dominance and allow the model to focus on the minority class.

##### 2. Deployment Preparation:

Following the optimization phase, a series of steps will ensure a successful deployment of the chosen model:

1. Once all testing and validation efforts are complete, the model configuration demonstrating the best performance will be selected for deployment.

Given the current model's seemingly perfect performance, exploring a simpler or pruned Random Forest variant holds promise.

- a) Simpler models often require fewer resources and compute power, leading to faster execution. This is crucial when considering deployment and scaling the model for real-world applications.
- b) While the current model achieves perfect scores, a simpler version might still deliver comparable performance, especially if the data is indeed very simple or dominated by a few strong features.

2. A comprehensive plan outlining the integration of the chosen model into the existing IT infrastructure should be formulated.

3. A monitoring system should be established to track the model's performance in a production environment using real-world data. This will be accompanied by a maintenance schedule for periodically updating the model to maintain its effectiveness.

4. Training materials and comprehensive documentation should be developed to educate end-users on the model's functionality, its intended use within business processes, and its limitations.

5. The deployment process should adhere to all relevant compliance regulations and data privacy rules, ensuring the integrity and security of customer data.

By prioritizing model optimization followed by a well-defined deployment plan, we aim to maximize the predictive power of the customer repurchase prediction model while guaranteeing a smooth and effective transition to real-world applications.



