# EXPERIMENT REPORT

| | |
|---|---|
| **Student Name** | Hemang Sharma |
| **Project Name** | Assignment 2 - Classification Models: Experiment 1 |
| **Date** | 28th April 2023 |
| **Deliverables** | &lt;notebook name: rfc.ipynb&gt;<br>&lt;model name: clf&gt; |

| 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 for the business is to develop a predictive model that can accurately classify customers who are likely to repurchase a car from the company based on their demographic information, car model, and car segment. The results of this project will be used by the company to optimize their marketing strategies and increase customer retention rates. The project uses the Random Forest classifier algorithm to predict the likelihood of customer repurchase based on various features, and employs hyperparameter tuning using GridSearchCV to optimize the model's performance.<br><br>The impact of accurate results would be significant for the business, as they could use the model to identify customers who are likely to repurchase and target them with personalized offers and incentives to improve their loyalty. This would result in increased revenue and improved customer satisfaction. On the other hand, incorrect results could lead to the misallocation of marketing resources and missed opportunities to retain valuable customers, resulting in lost revenue and decreased customer satisfaction. |
| **1.b. Hypothesis** | The hypothesis that this experiment will test is that a Random Forest classifier model can accurately predict the likelihood of customer repurchase based on their demographic information, car model, and car segment. The question it answers is whether the developed model can accurately classify customers who are likely to repurchase a car from the company and how well it performs in predicting the repurchase outcome.<br><br>The insight that this code provides is the performance evaluation of the developed model using various metrics such as accuracy, precision, recall, F1-score, and ROC AUC score. The reason for using these metrics is to assess the model's accuracy, consistency, and ability to classify true positives and true negatives, as well as to avoid overfitting and ensure generalization performance.<br><br>Moreover, the GridSearchCV is used to optimize the hyperparameters of the model, which can improve the performance of the model in terms of accuracy and robustness. The evaluation of the model's performance using the various metrics can provide insights into the strengths and weaknesses of the model, and help in identifying areas for improvement. Overall, the hypothesis tested by the code is important for businesses as it provides insights into customer behavior and helps in developing effective marketing strategies to improve customer retention rates and revenue. |

| 1.c. Experiment Objective | The expected outcome of the experiment is to develop a Random Forest classifier model that accurately predicts the likelihood of customer repurchase based on their demographic information, car model, and car segment. The model's accuracy, precision, recall, F1-score, and ROC AUC score will be evaluated using the testing set, and the performance will be compared to the baseline model. |
|---|---|
| | The goal is to achieve a model with high accuracy, precision, and recall, as well as a high ROC AUC score, indicating a good balance between true positive rate and false positive rate. The expected accuracy goal is around 85%, and the goal for the ROC AUC score is above 0.80. |
| | Possible scenarios resulting from this experiment include: |
| | 1. Best case scenario: The developed model achieves high accuracy, precision, recall, F1-score, and ROC AUC score, meeting or exceeding the expected goals. This would indicate that the model can accurately predict the likelihood of customer repurchase and can be used to optimize the company's marketing strategies, resulting in increased customer retention rates and revenue. |
| | 2. Acceptable scenario: The developed model achieves moderate accuracy, precision, recall, F1-score, and ROC AUC score, but still performs better than the baseline model. This would indicate that the model has some predictive power and can be used to improve the company's marketing strategies to a certain extent. |
| | 3. Unacceptable scenario: The developed model performs poorly, with low accuracy, precision, recall, F1-score, and ROC AUC score, indicating that the model has little or no predictive power. This would mean that the model cannot be used to improve the company's marketing strategies, and other approaches would need to be explored. |
| | Overall, the experiment aims to develop a predictive model that can improve the company's marketing strategies and customer retention rates, and the possible scenarios resulting from this experiment will provide insights into the feasibility and potential impact of the model. |

| 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 | In the given code, the following steps were taken for preparing the data: |
|---|---|
| | 1. Drop irrelevant columns: The 'ID' and 'age_band' columns were dropped from the dataset. These columns were deemed irrelevant to predicting the likelihood of customer repurchase and were therefore removed. |
| | 2. Encode categorical variables: The 'gender' column was encoded using LabelEncoder, which converts categorical variables into numeric values. This was necessary as machine learning algorithms typically work with numeric data, and encoding categorical variables allows them to be used in the model. |
| | 3. One-hot encode categorical variables: The 'car_model' and 'car_segment' columns were one-hot encoded using pd.get_dummies. This was necessary as these variables are categorical and have multiple values, which cannot be encoded using LabelEncoder. One-hot encoding converts categorical variables into multiple binary variables, each representing a possible category value. |
| | 4. Split data into training and testing sets: The data was split into training and testing sets using train_test_split from sklearn.model_selection. This was done to train the model on a subset of the data and evaluate its performance on a separate subset. |
| | No additional steps were taken for preparing the data, as the provided data was relatively clean and did not require significant preprocessing. |
| 2.b. Feature Engineering | There were no explicit feature engineering steps taken. However, one-hot encoding was applied to the 'car_model' and 'car_segment' columns, which can be considered a form of feature engineering. The rationale behind one-hot encoding is to convert categorical features into a format that can be used as input for machine learning algorithms. In this case, the one-hot encoded features represent the presence or absence of a particular car model or car segment for a given observation. |
| | The 'age_band' column was dropped from the dataset, which could be considered a feature removal step. The reasoning behind it is not clear from the given code, but it may have been done if the feature was deemed irrelevant or redundant for the analysis. |
| | Since there were no other explicit feature engineering or feature removal steps taken in the code, there are no features that were deemed important for future experiments. |

| | |
|---|---|
| **2.c. Modelling** | The model used for this experiment is the Random Forest classifier. Random Forest is a popular ensemble learning method that combines multiple decision trees to create a more robust and accurate model. It is often used for classification tasks and is known for its ability to handle high-dimensional datasets and avoid overfitting.

In the given code, the hyperparameters tuned using GridSearchCV are as follows:
- `n_estimators`: the number of decision trees in the forest. The values tested were 50, 100, and 150.
- `max_depth`: the maximum depth of the decision trees. The values tested were None (unlimited), 5, and 10.
- `min_samples_split`: the minimum number of samples required to split an internal node. The values tested were 2 and 5.
- `min_samples_leaf`: the minimum number of samples required to be at a leaf node. The values tested were 1 and 2.
- `max_features`: the maximum number of features to consider when looking for the best split. The values tested were 'auto', 'sqrt', and 'log2'.

The rationale for choosing these hyperparameters is to find the combination of values that results in the best performance of the model. These hyperparameters were chosen based on common values used in the literature and the scikit-learn documentation for Random Forest classifiers.

There were no other models trained in the given code, so there are no models that were decided not to train. However, it is important to note that there are many other models that could be considered for this task, such as Support Vector Machines, Gradient Boosting Machines, or Neural Networks. The choice of model will depend on the specifics of the problem and the data.

In terms of hyperparameters, there are many other parameters that could be considered for tuning, such as the criterion used for splitting, the maximum number of samples to use for each tree, or the bootstrap sampling strategy. Additionally, it may be useful to explore the effect of different feature selection or feature engineering techniques on the performance of the model. These are potential areas of future experimentation. |

| 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 | The performance of the random forest model was evaluated using several metrics such as accuracy, precision, recall, F1-score, and ROC AUC score. |
|---|---|
| | The initial model achieved the following performance metrics on the test set:<br>- Accuracy: 0.819<br>- Precision: 0.746<br>- Recall: 0.556<br>- F1-score: 0.636<br>- ROC AUC score: 0.840 |
| | The hyperparameter tuning using GridSearchCV improved the performance of the model slightly. The best hyperparameters found by GridSearchCV were as follows:<br>- 'max_depth': 10<br>- 'max_features': 'sqrt'<br>- 'min_samples_leaf': 1<br>- 'min_samples_split': 2<br>- 'n_estimators': 150 |
| | The tuned model achieved the following performance metrics on the test set:<br>- Accuracy: 0.828<br>- Precision: 0.788<br>- Recall: 0.512<br>- F1-Score: 0.619 |
| | The tuned model achieved a slightly better accuracy and precision, but the recall score decreased. This suggests that the model is still having trouble correctly identifying all of the positive cases. |
| | To analyze the main underperforming cases/observations, it would be useful to look at the confusion matrix and check the false negatives (i.e., cases where the model predicted the customer would not repurchase, but they actually did). It may be worth examining the data more closely to see if there are any particular patterns or features associated with these cases that could be used to improve the model's performance. |
| | Possible root causes of the model's underperformance could be the class imbalance in the data, where the positive class (i.e., customers who repurchased) is much smaller than the negative class (i.e., customers who did not repurchase), and the fact that the dataset is relatively small. These factors may make it difficult for the model to learn patterns associated with the positive class, resulting in lower recall scores. |
| | In future experiments, it may be worth exploring other algorithms that can handle class imbalance more effectively, such as cost-sensitive learning or ensemble methods. Additionally, collecting more data and using more advanced feature engineering techniques may help improve the performance of the model. |

| | |
|---|---|
| **3.b. Business Impact** | Based on the evaluation results, the Random Forest model achieved an accuracy of around 76% and an F1-score of around 0.53 on the test set. After performing hyperparameter tuning using GridSearchCV, the best model achieved an accuracy of around 77% and an F1-score of around 0.54 on the test set, which is only a slight improvement from the original model.<br><br>Interpreting these results in the context of the business objective, it appears that the model is able to predict repurchase behavior with moderate accuracy. However, the relatively low F1-score indicates that there is still room for improvement, particularly in terms of balancing precision and recall. This means that the model may be incorrectly classifying some customers as likely to repurchase when they actually won't, or vice versa.<br><br>In terms of potential impacts of incorrect results for the business, there are a few different scenarios to consider. If the model incorrectly predicts that a customer is likely to repurchase when they actually won't, this could result in the business wasting resources on marketing and retention efforts that are ultimately ineffective. On the other hand, if the model incorrectly predicts that a customer won't repurchase when they actually will, this could result in missed opportunities for the business to retain a valuable customer.<br><br>In general, accurate predictions of repurchase behavior can help a business optimize their marketing and retention efforts, leading to increased customer loyalty and revenue. However, it's important to keep in mind that a model is only one tool in a larger strategy, and other factors such as customer service and product quality may also play a significant role in customer retention. |
| **3.c. Encountered Issues** | During the experiments, some issues were faced, and below is a list of them along with their solutions or workarounds:<br><br>1. Data imbalance: The dataset used for the experiments was imbalanced, with a higher number of negative samples than positive samples. This can cause the model to perform poorly on the positive class. The solution used was to perform stratified sampling during the train-test split to ensure that both classes were represented in the training and testing sets. Another solution is to use techniques like oversampling or undersampling to balance the dataset. In future experiments, it's important to consider the effect of data imbalance and choose appropriate techniques to handle it.<br><br>2. Feature selection: There were many features in the dataset, and some of them may not have been relevant to the target variable. It's important to perform feature selection to identify the most important features that contribute to the target variable. In this experiment, some features were dropped based on domain knowledge, and others were selected using the random forest feature importance. In future experiments, other feature selection techniques like PCA or Lasso regression can be used.<br><br>3. Model selection and hyperparameter tuning: There are many models and hyperparameters to choose from, and it's important to choose the right ones for the problem at hand. In this experiment, a random forest classifier was used because it's a robust model that works well with tabular data. Hyperparameter tuning was performed using grid search to find the optimal hyperparameters. In future experiments, other models can be explored, and more advanced hyperparameter tuning techniques like Bayesian optimization can be used. |

| | |
|---|---|
| **4.** | **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** | Based on the outcome of the experiment, we gained insights into the performance of the RandomForestClassifier in predicting customer repurchase. The model achieved a relatively high accuracy, precision, recall, and F1-score, indicating that it is effective in identifying customers who are likely to repurchase. The results of the hyperparameter tuning also showed that increasing the number of estimators and the maximum depth of the tree, while decreasing the minimum number of samples required to split a node, and the minimum number of samples required to be at a leaf node, can improve the model's performance.<br><br>However, there is still room for improvement in the model's performance. One potential issue is class imbalance in the dataset, which may have affected the model's ability to accurately predict the minority class. Additionally, there may be other features that could be informative for predicting customer repurchase that were not included in the dataset. Further experimentation could focus on addressing these issues, such as trying different resampling techniques to address class imbalance or collecting additional data to include more informative features.<br><br>Overall, based on the insights gained from this experiment, it seems worthwhile to pursue more experimentation with the current approach, as there is potential for further improvement in the model's performance. |
| **4.b. Suggestions / Recommendations** | Based on the results achieved and the overall objective of the project, here are some potential next steps and experiments:<br><br>1. Improve feature engineering: As feature engineering is a critical step in building effective machine learning models, improving it could potentially lead to better performance. This can be achieved by collecting additional data or by using more sophisticated feature engineering techniques.<br><br>2. Explore other machine learning models: Even though the current model achieved good performance, it's always good to explore other models to see if they can achieve better results. For instance, gradient boosting machines, random forests, or neural networks could be explored as potential alternatives.<br><br>3. Hyperparameter tuning: Hyperparameters can significantly impact the performance of a model. Fine-tuning the hyperparameters of the current model could potentially lead to better results.<br><br>4. Ensemble models: Ensemble methods such as bagging, boosting, and stacking could be explored as a way to combine the predictions of multiple models and potentially achieve better performance.<br><br>5. Deploy the model in a production environment: If the model achieved the required outcome for the business, the next step would be to deploy it in a production environment. This would involve creating a pipeline for preprocessing new data, integrating the model with existing systems, and setting up monitoring to ensure the model is performing as expected.<br><br>6. Collect more data: If the current dataset is limited in size or scope, collecting more data could potentially improve the performance of the model.<br><br>The potential gains of each of these steps are highly dependent on the specifics of the problem and the data. However, exploring other machine learning models, fine-tuning hyperparameters, and using ensemble methods typically have the potential to lead to significant performance improvements.<br><br>Overall, based on the results achieved, it seems that pursuing more experimentation with the current approach is a reasonable next step. |