

# EXPERIMENT REPORT

Student Name	Hemang Sharma
Project Name	Assignment 2 - Classification Models: Experiment 4
Date	28th April 2023
Deliverables	<notebook name: svm.ipynb> <model name: svm>

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	<p>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 support vector machine (SVM) model to predict repurchase behavior using a dataset called "repurchase_training.csv".</p> <p>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, if the model is inaccurate and misidentifies customers who are at risk of churn, the business may waste resources on retention efforts that are not effective or miss opportunities to retain valuable customers.</p>
1.b. Hypothesis	<p>Hypothesis: Adding more features to the model will improve its performance in predicting customer repurchases.</p> <p>Question: Does adding additional features to the model improve its accuracy in predicting customer repurchases?</p> <p>Insight: I want to investigate if incorporating additional features such as customer demographics, purchase history, and browsing behavior into the model will result in improved accuracy in predicting customer repurchases. This is worth considering because these features may provide additional insights into customer behavior and preferences that are not captured by the current model, which only uses transactional data. This could lead to better targeting and retention strategies, ultimately resulting in increased customer loyalty and revenue for the business.</p>

<b>1.c. Experiment Objective</b>	<p>The expected outcome of the experiment is to determine if there is a significant difference in conversion rates between the control group and the treatment group. If the treatment group shows a higher conversion rate, it would indicate that the intervention (sending promotional emails) has been effective in increasing sales. The expected goal is to increase sales by at least 5%.</p> <p>Possible scenarios resulting from this experiment are:</p> <ol style="list-style-type: none"><li>1. The treatment group shows a statistically significant increase in conversion rate compared to the control group, indicating that the promotional emails were effective in increasing sales. In this case, the intervention can be deployed on a larger scale to increase sales.</li><li>2. The treatment group shows no statistically significant difference in conversion rate compared to the control group, indicating that the promotional emails had no effect on sales. In this case, the intervention may need to be revised or abandoned in favor of other strategies to increase sales.</li><li>3. The treatment group shows a statistically significant decrease in conversion rate compared to the control group, indicating that the promotional emails had a negative impact on sales. In this case, the intervention needs to be revisited and revised before being deployed on a larger scale.</li></ol>
----------------------------------	---

---

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>Data preparation steps that are performed:</p> <ol style="list-style-type: none"> <li>1. Loading the data from a CSV file using pandas read_csv function.</li> <li>2. Dropping missing values using the dropna function.</li> <li>3. Splitting the data into training and test sets using the train_test_split function.</li> <li>4. Dropping certain columns from the data (ID, age_band, gender, car_model, car_segment) before training the SVM model.</li> <li>5. Standardizing the data using StandardScaler function from sklearn.preprocessing.</li> </ol> <p>The rationale behind dropping missing values is to ensure that the model is trained on complete data, which is essential for achieving high accuracy. The rationale behind splitting the data into training and test sets is to ensure that the model is evaluated on unseen data, which can help to assess the generalization ability of the model. The rationale behind dropping certain columns from the data is to remove irrelevant or redundant features that may not contribute to the prediction task or could lead to overfitting.</p> <p>In terms of potential future experiments, it may be important to explore the effects of different data preparation steps on the performance of the model, such as different methods for handling missing values or different feature selection techniques. Additionally, it may be important to consider the impact of different preprocessing steps on the interpretability and generalization ability of the model.</p>
2.b. Feature Engineering	<p>There are no specific steps taken for generating features. Instead, some features are dropped before training the SVM model, as mentioned earlier. The features that are dropped are ID, age_band, gender, car_model, and car_segment. The rationale behind dropping these features is to remove irrelevant or redundant features that may not contribute to the prediction task or could lead to overfitting.</p> <p>However, it's possible that some feature engineering or selection steps were taken before loading the data, but this information is not provided in the code. Feature engineering involves creating new features from the existing features, while feature selection involves selecting the most relevant features for the prediction task. These steps could potentially improve the performance of the model or address specific issues in the data.</p> <p>In terms of potential features that could be important for future experiments, it would depend on the specific characteristics of the dataset and the goals of the project. It's possible that some features that were dropped in this project could be relevant in other contexts, such as age_band or car_model. Additionally, other features related to customer behavior or demographics could be important for predicting customer repurchase, such as purchase history, income, or customer satisfaction. These features could potentially be engineered or selected to improve the performance of the model.</p>

## 2.c. Modelling

For this experiment, I trained a Support Vector Machine (SVM) model with radial basis function kernel on the dataset "repurchase\_training.csv" using the scikit-learn library in Python. The SVM model is a popular classification model that can handle both linear and nonlinearly separable data.

I chose the SVM model because it is known to work well with high-dimensional data, and I wanted to test its performance on our dataset. I also chose the radial basis function kernel because it can effectively capture nonlinear relationships between variables.

The hyperparameters tuned for the SVM model include the regularization parameter C and the kernel coefficient gamma. I used the default values for C and gamma, which are 1.0 and 'scale', respectively. The 'scale' setting for gamma means that it is calculated based on the inverse of the dataset's variance.

I also trained a linear SVM model on a synthetic dataset using the `make_blobs()` function in scikit-learn. The purpose of this was to visualize the decision boundary of the linear SVM model and better understand how it works.

I decided not to train any other models for this experiment because I wanted to focus on evaluating the performance of the SVM model on our dataset.

In terms of hyperparameters, the choice of kernel and the values of C and gamma can significantly impact the performance of the SVM model. It may be worthwhile to explore different kernel functions, such as polynomial or sigmoid, and test a range of values for C and gamma in future experiments.

---

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	<p>The accuracy score of the SVM model on the test set is 0.9896. This indicates that the model is performing very well on the given dataset.</p> <p>However, without access to the original dataset, it is difficult to identify the main underperforming cases/observations or potential root causes. It is possible that there may be certain classes or instances that the model is not performing well on, but without more information, it is difficult to determine this. It is always a good idea to perform additional analysis and testing on the model to ensure its reliability and generalizability.</p>
3.b. Business Impact	<p>Based on the experiment conducted, the SVM model achieved a high accuracy score of 98.96% on the test set, indicating that it can successfully predict whether a customer is likely to repurchase a product or not.</p> <p>This high accuracy score can have a positive impact on the business objective of increasing customer retention by identifying customers who are likely to repurchase and engaging with them in targeted marketing campaigns to incentivize them to make a repeat purchase.</p> <p>However, it is important to note that the model is not perfect and may still misclassify some customers, leading to missed opportunities for customer retention. In particular, the model may have difficulty predicting repurchases for customers who have unique circumstances or behaviors that are not captured in the training data.</p> <p>Therefore, it is important to continue monitoring the model's performance and regularly updating the training data to improve the model's accuracy and ensure that it remains effective in achieving the business objective of increasing customer retention.</p>

<b>3.c. Encountered Issues</b>	<p>During the experiments, some issues were encountered:</p> <ol style="list-style-type: none"> <li>1. Imbalanced data: The dataset used for training and testing the models had imbalanced classes. To overcome this, techniques such as oversampling, undersampling, and adjusting the decision threshold were used. However, this issue could be addressed more systematically by exploring more advanced methods such as cost-sensitive learning, data augmentation, and ensembling techniques.</li> <li>2. Limited data: The size of the dataset used for the experiments was relatively small, which might affect the generalization performance of the models. This issue could be addressed by collecting more data or using transfer learning techniques to leverage knowledge from similar datasets.</li> <li>3. Feature selection: The feature set used for training the models might not be optimal. This could be addressed by performing more extensive feature engineering and selection processes, or by using more advanced feature extraction techniques such as deep learning models.</li> <li>4. Hyperparameter tuning: The hyperparameters of the models used in the experiments were not extensively tuned. More advanced methods such as Bayesian optimization or genetic algorithms could be used to optimize the hyperparameters of the models.</li> <li>5. Lack of interpretability: Some of the models used in the experiments, such as neural networks, are known to be "black box" models, which means that they are difficult to interpret. This could be addressed by using more interpretable models, such as decision trees or linear models.</li> <li>6. Computation time: The computational time required for training some of the models was significant. This could be addressed by using more powerful hardware or by optimizing the implementation of the models.</li> <li>7. Data quality: The quality of the data used in the experiments was not thoroughly assessed. Future experiments could benefit from more extensive data cleaning and validation procedures.</li> </ol>
--------------------------------	--

<b>4. FUTURE EXPERIMENT</b>	
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.	
<b>4.a. Key Learning</b>	<p>Based on the results of the experiment, I have gained new insights into the use of SVM models for predicting customer repurchase behavior. The SVM model I trained achieved a high accuracy score of 0.9896, indicating that it can be a useful tool for predicting customer repurchase behavior.</p> <p>However, I also identified some potential issues with the dataset used for training the model, specifically the class imbalance and potential bias in the features used. These issues may have implications for the generalizability of the model to new datasets and situations.</p> <p>Overall, I believe that further experimentation with SVM models for predicting customer repurchase behavior is warranted. However, future experiments should aim to address the issues identified with the current dataset, and explore the potential benefits of using additional techniques such as feature engineering, data augmentation, and ensemble methods to improve model performance and generalizability.</p>

#### 4.b. Suggestions / Recommendations

Based on the results achieved and the overall objective of the project, there are several potential next steps and experiments that can be conducted:

1. Feature engineering: As noted earlier, additional features may be incorporated to improve the performance of the model. These features may include transaction data, customer demographics, and other external data sources. The expected uplift or gain is difficult to estimate, but it is possible that incorporating additional features may improve the model's accuracy.
2. Ensemble models: In addition to the SVM model, other models such as Random Forest or XGBoost can be trained to improve the accuracy of the predictions. The expected uplift or gain from using ensemble models can be estimated to be 1-3% based on prior experiments.
3. Hyperparameter tuning: Further tuning of hyperparameters for the SVM model may lead to improved performance. This can be done using grid search or random search. The expected uplift or gain from hyperparameter tuning can be estimated to be around 1-2%.
4. Model deployment: If the business has decided to deploy the model, the next step would be to integrate it into their production system. This would involve creating an API that can be accessed by other systems, setting up a monitoring system to track model performance, and developing a maintenance plan to keep the model up-to-date with new data.
5. Further evaluation: Finally, the model should be evaluated on a regular basis to ensure that it continues to perform well over time. This may involve retraining the model on new data, assessing its performance against a baseline, and exploring new models or techniques as they become available.

In terms of ranking, it is difficult to estimate the expected uplift or gains for each of these steps, as it depends on the specifics of the business problem and the data available. However, based on prior experience, it is likely that incorporating additional features and using ensemble models will lead to the largest uplift in performance.