

# Predicting customer lifetime value using the online retail

## **Problem Description :**

The objective of this project is to predict the customer lifetime value for an online retail business. Customer lifetime value is a metric that helps businesses understand the total value a customer will bring to their business over the entire relationship with that customer. By predicting customer lifetime value, businesses can make informed decisions about how much to invest in acquiring new customers and retaining existing customers.

The project will use the Online Retail dataset from Kaggle, which contains transaction data for a UK-based online retail company from 2010 to 2011. The dataset includes information such as customer ID, product description, quantity, price, and transaction date. You can find the dataset on Kaggle using this link: <https://www.kaggle.com/vijayuv/onlineretail>

## **Background Information**

Customer lifetime value is an important metric for any business, but it is particularly important for online retail businesses. These businesses often have lower profit margins and higher customer acquisition costs than traditional retail businesses. By predicting customer lifetime value, online retail businesses can optimize their marketing and retention strategies to maximize their profits and minimize their costs.

# **General Framework and Steps**

To solve this problem, we can follow these general steps:

1. **Data Cleaning and Preprocessing:** We will clean and preprocess the Online Retail dataset to ensure that the data is ready for analysis. This will involve handling missing values, removing duplicates, and transforming the data as needed.
2. **Feature Engineering:** We will create new features from the existing data that may be relevant to predicting customer lifetime value. This may include features such as total spend, number of purchases, and recency of purchases.
3. **Model Selection:** We will select a regression model that is suitable for predicting customer lifetime value. This may include models such as linear regression, decision tree regression, or random forest regression.
4. **Model Training and Evaluation:** We will train the selected model on a training dataset and evaluate its performance using a validation dataset. We will use metrics such as mean squared error and R-squared to evaluate the model's performance.
5. **Hyperparameter Tuning:** We will tune the hyperparameters of the selected model to optimize its performance.
6. **Model Deployment:** We will deploy the selected model in a production environment, where it can be used to predict customer lifetime value for new customers.

## **Deliverables**

The deliverables for this project will include:

1. A cleaned and preprocessed version of the Online Retail dataset
2. A set of relevant features for predicting customer lifetime value
3. A trained regression model that accurately predicts customer lifetime value
4. Documentation of the model selection, training, and evaluation process
5. A deployed model that can be used to predict customer lifetime value for new customers

## **Code Explanation :**

Here is the simple explanation for the code which is provided in the code.py file.

### **Section 1: Data Cleaning and Preprocessing**

In this section, we load the Online Retail dataset and clean and preprocess the data. We remove rows with missing values and duplicates, convert the InvoiceDate column to a datetime object, and remove rows with negative quantity or price. We also create new columns for total spend and recency of purchases, and group the data by customer ID and aggregate the TotalSpend and Recency columns. The cleaned and preprocessed data is then saved to a new variable called df.

### **Section 2: Feature Engineering and Model Selection**

In this section, we create new features from the existing data that may be relevant to predicting customer lifetime value. We create a new column for frequency of purchases, a new column for average spend per transaction, and a new column for customer lifetime value. We then select relevant columns for training the model, split the data into training and validation sets, and train a linear regression model, a decision tree regression model, and a random forest regression model. We evaluate the models on the validation set and print the evaluation metrics for each model.

### **Section 3: Hyperparameter Tuning and Model Deployment**

In this section, we tune the hyperparameters of the selected model to optimize its performance. We perform a grid search to find the best hyperparameters for the random forest regression model and print the best hyperparameters found by grid search. We then train the final random forest regression model with the best hyperparameters and save the model to a file called final\_model.pkl.

## **Future Work :**

To further improve the accuracy of the customer lifetime value prediction model, there are several additional steps that can be taken:

### **Step 1: Feature Selection**

In this step, we can use feature selection techniques to identify the most important features for predicting customer lifetime value. This will allow us to focus on the most relevant features and exclude the ones that are not useful. Feature selection techniques such as correlation analysis, mutual information, and recursive feature elimination can be used to perform feature selection.

### **Step 2: Model Ensemble**

In this step, we can combine multiple models to create a model ensemble that is more accurate than any individual model. Model ensemble techniques such as bagging, boosting, and stacking can be used to create a model ensemble. Bagging involves training multiple models on different subsets of the data and averaging their predictions, while boosting involves training multiple models sequentially and adjusting the weights of misclassified samples. Stacking involves combining the predictions of multiple models using a meta-learner.

### **Step 3: Time Series Analysis**

In this step, we can use time series analysis techniques to model the trends and seasonality in the data. This will allow us to account for the temporal nature of the data and make more accurate predictions. Time series analysis techniques such as ARIMA, SARIMA, and Prophet can be used to model the time series data.

### **How to Implement the Future Work**

**To implement the future work, we can follow the following step-by-step guide:**

1. Perform feature selection on the existing features to identify the most important ones for predicting customer lifetime value.
2. Train multiple models using the selected features and evaluate their performance on the validation set.
3. Use model ensemble techniques to combine the predictions of the multiple models and create a more accurate model.
4. Perform time series analysis on the data to model the trends and seasonality in the data.
5. Combine the time series analysis with the existing models and create a final model that is more accurate than the existing model.

By following this step-by-step guide, we can improve the accuracy of the customer lifetime value prediction model and make more informed decisions about how to invest in acquiring and retaining customers.

## **Exercise :**

**Try to answers the following questions by yourself to check your understanding for this project. If stuck, detailed answers for the questions are also provided.**

### **1.What is customer lifetime value and why is it important for businesses?**

Answer: Customer lifetime value (CLV) is a metric that helps businesses understand the total amount of revenue they can expect to receive from a single customer over the course of their relationship. It takes into account the frequency of purchases, the average spend per transaction, and the recency of purchases. CLV is important for businesses because it helps them make informed decisions about how much to invest in acquiring new customers and retaining existing customers. By predicting the CLV of a customer, a business can decide how much they are willing to spend to acquire or retain that customer.

### **2.What is the purpose of feature engineering in this project and what features are created?**

Answer:The purpose of feature engineering in this project is to create new features from the existing data that may be relevant to predicting customer lifetime value. The features created in this project are the frequency of purchases, the average spend per transaction, and the customer lifetime value. The frequency of purchases is calculated by counting the number of transactions a customer has made. The average spend per transaction is calculated by dividing the total spend by the frequency of purchases. The customer lifetime value is calculated by multiplying the average spend per transaction by the frequency of purchases and dividing by the recency of purchases.

### **3.What is the difference between a decision tree and a random forest?**

Answer: A decision tree is a machine learning algorithm that creates a tree-like model of decisions and their possible consequences. Each internal node of the tree represents a decision, and each leaf node represents a possible outcome. A random forest is an ensemble learning method that creates a large number of decision trees and combines their predictions. The idea

behind a random forest is that by creating multiple decision trees and combining their predictions, we can reduce the variance and improve the accuracy of the model.

#### **4.What is hyperparameter tuning and why is it important?**

Answer: Hyperparameter tuning is the process of selecting the best hyperparameters for a machine learning model. Hyperparameters are values that are set before training the model and are not learned from the data. Examples of hyperparameters include the number of trees in a random forest, the learning rate in a gradient boosting model, and the regularization parameter in a linear regression model. Hyperparameter tuning is important because it can significantly improve the performance of a machine learning model. By selecting the best hyperparameters, we can create a model that is more accurate and more robust to new data.

#### **5.What is model deployment and why is it important?**

Answer: Model deployment is the process of making a trained machine learning model available for use in a production environment. This involves taking the trained model and integrating it into a larger software system, such as a web application or a mobile app. Model deployment is important because it allows businesses to use the trained model to make predictions on new data and make informed decisions based on those predictions. Without model deployment, the trained model would only be useful for offline analysis and would not be able to provide real-time predictions on new data.