# Predicting Customer lifetime value

## Problem Statement

**Background:** Customer Lifetime Value (CLV) is a crucial metric for businesses to understand the long-term value of their customers. It represents the predicted net profit generated by a customer over their entire relationship with the company. By accurately predicting CLV, businesses can make informed decisions regarding customer acquisition, retention, and resource allocation.

The objective of this project is to develop a predictive model that can estimate the CLV of customers. By leveraging customer data and applying machine learning techniques, the model will provide insights into the potential value of individual customers and enable businesses to optimize their marketing, sales, and customer relationship management strategies.

**Dataset:** The dataset used for this project is "customer_segmentation.csv". It contains transactional data from an e-commerce company, including information about customer purchases. The dataset includes the following columns:

- CustomerID: Unique identifier for each customer.
- InvoiceNo: Unique identifier for each invoice.
- InvoiceDate: Date and time of the invoice.
- UnitPrice: Price per unit of the purchased item.
- Quantity: Quantity of items purchased.
- Revenue: Total revenue generated by the purchase.
- Country: Country where the purchase was made.
- Methodology: The project follows the following steps:

**Data Preprocessing:** The dataset is loaded and preprocessed to convert data types, extract relevant features, and create additional fields for analysis.

1. Recency Calculation: The recency of each customer is calculated by determining the number of days since their last purchase. K-means clustering is then applied to assign a recency score to each customer.
2. Frequency Calculation: The total number of orders made by each customer is calculated to determine their frequency. K-means clustering is used to assign a frequency score to each customer.
3. Revenue Calculation: The total revenue generated by each customer is calculated. K-means clustering is applied to assign a revenue score to each customer.
4. Overall Score Calculation: The recency, frequency, and revenue scores are combined to calculate an overall score for each customer, representing their value to the business.
5. Customer Segmentation: Customers are segmented into three categories based on their overall scores: Low-Value, Mid-Value, and High-Value.
6. Customer Lifetime Value Prediction: A subset of the dataset is used to calculate the 6-month revenue for each customer. Outliers and negative values are filtered out, and the remaining data is merged with the customer segmentation data. Correlations between CLV and the feature set are analyzed.

**Model Development:** Machine learning models, such as regression or XGBoost, are trained using the feature set and CLV values. The models aim to predict the CLV of customers based on their characteristics and historical data.

**Model Evaluation:** The trained models are evaluated using appropriate metrics, such as mean squared error or R-squared, to assess their predictive performance. The best-performing model is selected for CLV prediction.

**CLV Prediction:** The selected model is used to predict the CLV of new or existing customers. These predictions can be used to guide business decisions related to customer segmentation, marketing strategies, and resource allocation.

**Conclusion:** Predicting Customer Lifetime Value is essential for businesses to effectively allocate resources, optimize marketing strategies, and enhance customer relationship management. By leveraging customer data and employing machine learning techniques, businesses can gain valuable insights into the long-term value of their customers. This project aims to develop a predictive model that estimates CLV and provides actionable insights for business decision-making.

Note: The code provided in the previous conversation demonstrates the implementation of the steps outlined in the methodology section.

# Framework

**Import Necessary Libraries:** Begin by importing the required libraries for data manipulation, analysis, and machine learning, such as pandas, numpy, scikit-learn, and XGBoost.

**Load and Preprocess the Dataset:** Load the "customer_segmentation.csv" dataset using pandas' read_csv() function. Perform data preprocessing steps like converting data types, handling missing values, and creating additional features if necessary. Ensure the dataset is in a suitable format for further analysis.

**Recency Calculation:** Calculate the recency of each customer by determining the number of days since their last purchase. Use pandas' groupby() function to group the data by CustomerID and find the maximum InvoiceDate for each customer. Then, calculate the difference between the maximum InvoiceDate and the current date to obtain the recency value for each customer.

**Frequency Calculation**: Calculate the total number of orders made by each customer to determine their frequency. Again, use pandas' groupby() function to group the data by CustomerID and count the unique InvoiceNo values for each customer.

**Revenue Calculation:** Calculate the total revenue generated by each customer. Group the data by CustomerID using groupby() and sum the Revenue column for each customer.

**K-means Clustering for Recency, Frequency, and Revenue:** Apply K-means clustering to assign scores to customers based on their recency, frequency, and revenue values. Use scikit-learn's KMeans class to perform clustering separately for each attribute. Choose an appropriate number of clusters and fit the data. Assign the cluster labels to each customer.

**Overall Score Calculation:** Combine the recency, frequency, and revenue scores to calculate an overall score for each customer. You can assign weights to each score based on business requirements and calculate the weighted sum to obtain the overall score.

**Customer Segmentation:** Segment customers into three categories based on their overall scores. For example, define thresholds to classify customers as Low-Value, Mid-

Value, and High-Value. Create a new column in the dataset to store the customer segmentation labels.

**Subset and Prepare Data for CLV Prediction:** Select a subset of the dataset to calculate the 6-month revenue for each customer. Filter out outliers and negative values, as they can skew the CLV predictions. Merge this data with the customer segmentation data obtained in the previous step.

**Model Development:** Split the data into training and testing sets. Choose a suitable machine learning algorithm for CLV prediction, such as linear regression or XGBoost. Train the model using the training data, considering the feature set as input and the CLV values as the target variable.

**Model Evaluation:** Evaluate the trained model using appropriate evaluation metrics like mean squared error (MSE), root mean squared error (RMSE), or R-squared. Calculate the predictions for the testing data and compare them with the actual CLV values to assess the model's performance.

**CLV Prediction:** Once the model is trained and evaluated, it can be used to predict the CLV of new or existing customers. Provide the necessary input data, such as customer characteristics and historical data, to the trained model, and obtain the predicted CLV values.

**Further Analysis and Interpretation:** Perform further analysis on the predicted CLV values, such as calculating summary statistics, visualizing the distribution of CLV across customer segments, and exploring correlations between CLV and other features. This analysis can provide valuable insights for business decision-making.

**Conclusion and Documentation:** Summarize the findings and conclusions of the CLV prediction project. Document the code, methodology, and results for future reference.

# Code Explanation

**Section 1:** Importing Libraries In this section, the necessary libraries for data manipulation and machine learning are imported. These libraries provide functions and tools to perform various tasks, such as handling data, training machine learning models, and evaluating their performance.

**Section 2:** Loading and Preprocessing the Dataset This section involves loading the dataset from the "customer_segmentation.csv" file into a pandas DataFrame. The dataset contains information about customers, their purchases, and revenue generated. The read_csv() function from the pandas library is used to load the data.

Once the data is loaded, it goes through preprocessing steps. The InvoiceDate column is converted to a datetime data type for easier manipulation. Missing values, if any, are handled, and additional features can be created if required. These steps ensure that the data is in a suitable format for analysis.

**Section 3:** Recency Calculation Recency refers to the number of days since a customer's last purchase. In this section, the code calculates the recency value for each customer by finding the maximum InvoiceDate for each unique CustomerID using pandas' groupby() function. The difference between the maximum InvoiceDate and the current date is then calculated to obtain the recency value.

**Section 4:** Frequency Calculation Frequency indicates the total number of orders made by each customer. The code calculates the frequency by counting the unique InvoiceNo values for each customer using pandas' groupby() and nunique() functions. This step provides insights into how often a customer makes purchases.

**Section 5:** Revenue Calculation Revenue represents the total amount of money spent by each customer. The code calculates the revenue by summing the Revenue column for each customer using pandas' groupby() and sum() functions. This information helps understand the value generated by each customer.

**Section 6:** K-means Clustering for Recency, Frequency, and Revenue K-means clustering is a technique that groups similar data points together. In this section, the code applies K-means clustering separately for the recency, frequency, and revenue attributes. The KMeans class from the scikit-learn library is used to perform clustering.

The number of clusters can be chosen based on business requirements. The data is fitted to the K-means model, and cluster labels are assigned to each customer. Clustering helps identify different segments or groups of customers based on their recency, frequency, and revenue values.

**Section 7:** Overall Score Calculation In this section, the code calculates an overall score for each customer by combining the recency, frequency, and revenue scores. The scores can be weighted based on their importance to the business. By calculating the weighted sum of the scores, an overall score is obtained, which provides a comprehensive measure of a customer's value.

**Section 8:** Customer Segmentation Segmenting customers into different categories helps in understanding their value and tailoring marketing strategies accordingly. In this section, the code classifies customers into segments based on their overall scores. Thresholds can be defined to classify customers as Low-Value, Mid-Value, and High-Value. A new column is created in the dataset to store the customer segmentation labels.

**Section 9:** Subset and Prepare Data for CLV Prediction To predict the Customer Lifetime Value (CLV), a subset of the dataset is selected. This subset contains relevant features and the target variable, which is the 6-month revenue for each customer. Outliers and negative values are filtered out, as they can adversely affect the CLV predictions. The dataset is merged with the customer segmentation information to have a comprehensive view.

**Section 10:** Train/Test Split and Model Training To evaluate the performance of the CLV prediction model, the dataset is split into a training set and a testing set. The training set is used to train the machine learning model, while the testing set is used to assess its predictive capability. The data is divided into input features (X) and the target variable (y).

In this code, a Linear Regression model from the scikit-learn library is chosen for CLV prediction. The model is fitted to the training data, learning the patterns and relationships between the input features and the target variable.

**Section 11:** Model Evaluation Once the model is trained, it is evaluated using the testing set. The model predicts the CLV values for the testing set, and various evaluation metrics, such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared

score, are calculated. These metrics provide insights into the accuracy and performance of the CLV prediction model.

**Section 12:** Conclusion This code provides a step-by-step workflow to analyze customer data, segment customers based on their recency, frequency, and revenue, and predict their CLV using a machine learning model. The analysis helps identify high-value customers, tailor marketing strategies, and make data-driven business decisions.

# Future Work

Customer Lifetime Value (CLV) prediction and segmentation provide valuable insights into customer behavior and help businesses make data-driven decisions. Here are some future steps to enhance the project:

**1. Feature Engineering:**

- Explore additional customer features: Consider incorporating more customer attributes such as demographics, purchase history, customer interactions, or website engagement data. These features can provide deeper insights into customer behavior and improve CLV prediction accuracy.
- Create temporal features: Develop time-based features such as average purchase frequency, seasonal patterns, or customer tenure. These features can capture temporal variations in customer behavior and enhance CLV predictions.

**2. Advanced CLV Models:**

- Experiment with advanced regression models: Apart from Linear Regression, try other regression models such as Random Forest, Gradient Boosting, or Neural Networks. These models can capture complex relationships and improve CLV predictions.
- Explore customer lifetime models: Investigate advanced CLV models like Pareto/NBD or Gamma-Gamma models that estimate CLV based on customer purchase patterns and monetary value. These models can provide more accurate and granular predictions.

**3. Model Optimization:**

- Hyperparameter tuning: Fine-tune the hyperparameters of the chosen regression model to optimize its performance. Use techniques like Grid Search or Randomized Search to find the best combination of hyperparameters.
- Feature selection: Conduct feature selection techniques, such as Recursive Feature Elimination (RFE) or L1 regularization, to identify the most relevant features for CLV prediction. This can improve model performance and reduce overfitting.

**4. Customer Segmentation Enhancement:**

- Refine segmentation criteria: Evaluate the current segmentation approach and consider incorporating additional attributes or clustering algorithms (e.g., DBSCAN, Hierarchical Clustering) to create more meaningful customer segments. This can lead to more targeted marketing strategies.
- Conduct segmentation validation: Validate the effectiveness of the customer segmentation by analyzing the differences in customer behavior, purchasing patterns, or revenue across different segments. This validation will help ensure the segmentation accurately reflects customer characteristics.

**5. CLV Prediction Model Evaluation:**

- Cross-validation: Perform cross-validation to assess the stability and robustness of the CLV prediction model. Split the dataset into multiple folds and evaluate the model's performance on each fold to obtain more reliable metrics.
- Compare multiple models: Compare the performance of different regression models using various evaluation metrics (e.g., MAE, MSE, R-squared). This analysis will help choose the most suitable model for CLV prediction.

**Step-by-Step Guide to Implement the Future Work:**

- Data Exploration and Feature Analysis:
- Identify additional customer features that can provide valuable insights into customer behavior and enhance CLV predictions.
- Analyze the data to understand temporal patterns and create relevant time-based features.
- Ensure the new features are relevant, accurate, and representative of customer behavior.
- Advanced CLV Model Implementation:
- Select and implement advanced regression models such as Random Forest, Gradient Boosting, or Neural Networks.
- Adapt the code to incorporate the chosen model and evaluate its performance on CLV prediction.
- Compare the results with the existing Linear Regression model to assess improvement.

**Model Optimization:**

- Perform hyperparameter tuning for the chosen advanced model using techniques like Grid Search or Randomized Search.
- Implement feature selection techniques such as Recursive Feature Elimination or L1 regularization to identify the most relevant features for CLV prediction.
- Evaluate the optimized model's performance and compare it with the previous version.

**Refine Customer Segmentation:**

- Review the current customer segmentation criteria and identify opportunities for improvement.
- Incorporate additional customer attributes or explore different clustering algorithms to create more meaningful customer segments.
- Validate the effectiveness of the refined segmentation by analyzing customer behavior and characteristics within each segment.

**Evaluate CLV Prediction Models:**

- Perform cross-validation on the CLV prediction model to assess its stability and generalization capability.
- Calculate various evaluation metrics, such as MAE, MSE, and R-squared, for each model to compare their performance.
- Select the most suitable CLV prediction model based on the evaluation results.

By following this step-by-step guide, you can enhance the project by incorporating additional features, implementing advanced CLV models, optimizing the models, refining customer segmentation, and evaluating the performance of the CLV prediction models. These improvements will provide more accurate and insightful predictions, enabling businesses to make informed decisions and strategies based on customer behavior and value.

# Concept Explanation

Alright, let's dive into the magical world of Linear Regression! Imagine you're a wizard who wants to predict the price of a wand based on its length. How would you do it? Linear Regression comes to your rescue!

**The Magical Wand Shop Dataset:** You have collected data on various wands from your magical wand shop. Each wand has a length and a corresponding price. We want to find a way to predict the price of a new wand just by knowing its length. Let's cast our spell of Linear Regression!

**The Idea Behind Linear Regression:** Linear Regression is like drawing a straight line through a scatterplot of points. The line represents the relationship between the wand length and its price. It's like finding the best-fitting magic spell that connects the length and price of the wands.

The Wand Shop Adventure: Imagine you have a scatterplot of wand lengths and prices in front of you. You grab your magic wand (no pun intended) and draw a line that seems to pass close to most of the points. You want this line to be the best representation of the relationship between wand length and price.

**The Magical Equation:** $y = mx + b$ In the world of Linear Regression, the equation $y = mx + b$ holds great power. Here's what it means:

- $y$ is the predicted price of a wand.
- $m$ is the slope of the line. It determines how steep the line is. A positive slope means that as the length increases, the price increases too.
- $x$ is the length of the wand.
- $b$ is the y-intercept of the line. It tells us where the line intersects the y-axis. In our case, it represents the base price of a wand with a length of 0.

**Calculating the Magic Slope and Intercept:** To find the best values for $m$ and $b$, we unleash some mathematical wizardry. We use a technique called "Ordinary Least Squares" to minimize the distance between each data point and the line. The wizard's goal is to make the line as close as possible to all the points.

- Prediction Time: Once we have our magical values for $m$ and $b$, we can predict the price of a new wand just by knowing its length. You simply plug in the length

into the equation y = mx + b, and voila! The magic of Linear Regression gives you the predicted price.

- Example Incantation: Let's say the equation we found is y = 10x + 5. If you have a wand with a length of 7, you can predict its price like this:
- Plug in x = 7 into the equation: y = 10 * 7 + 5
- Calculate: y = 70 + 5

Abracadabra! The predicted price for a wand with length 7 is 75!

**The Power of Linear Regression:** Linear Regression helps us understand the relationship between two variables and make predictions based on that relationship. It's like having a crystal ball that can estimate prices, predict trends, or even foresee the future! It's a valuable tool in the magical world of data analysis.

So, embrace the power of Linear Regression and let it guide you in unraveling the mysteries hidden within your data. May your wand always be accurate and your predictions enchantingly precise!

# Exercise Questions

**Question: What is Linear Regression, and how does it work?**

**Answer:** Linear Regression is a statistical algorithm used to model the relationship between a dependent variable (target variable) and one or more independent variables (predictor variables). It assumes a linear relationship between the variables and aims to find the best-fitting line that minimizes the distance between the data points and the line. The equation of a simple linear regression model is y = mx + b, where y is the dependent variable, x is the independent variable, m is the slope, and b is the y-intercept.

**Exercise 2:** Model Evaluation

**Question: How do you evaluate the performance of a Linear Regression model?**

**Answer:** To evaluate the performance of a Linear Regression model, several metrics can be used:

1. Mean Squared Error (MSE): It calculates the average squared difference between the predicted and actual values. A lower MSE indicates a better fit.
2. Root Mean Squared Error (RMSE): It is the square root of the MSE, providing a measure in the original scale of the target variable.
3. R-squared ($R^2$) score: It represents the proportion of the variance in the dependent variable that can be explained by the independent variables. Higher values (closer to 1) indicate a better fit.
4. Adjusted R-squared score: It adjusts the R-squared score by the number of predictors in the model, providing a more reliable measure of model performance.

**Exercise 3:** Handling Categorical Variables

**Question: How can you include categorical variables in a Linear Regression model?**

**Answer:** Including categorical variables in a Linear Regression model requires encoding them into numerical form. One common technique is one-hot encoding, where each category is represented by a binary column. For example, if we have a categorical variable "Color" with categories red, blue, and green, we create three binary columns (red, blue, green). If an observation belongs to the red category, the red column is set to

1, and the others are set to 0. These new columns can then be used as independent variables in the regression model.

**Exercise 4:** Dealing with Overfitting

**Question: What is overfitting in the context of Linear Regression, and how can you address it?**

**Answer:** Overfitting occurs when a model learns the training data too well and performs poorly on unseen data. In Linear Regression, it can happen when the model becomes too complex and captures noise or irrelevant patterns. To address overfitting, we can:

1. Regularization: Introduce a penalty term to the regression equation, such as Ridge or Lasso regression, which restricts the magnitude of the coefficients.
2. Feature Selection: Remove irrelevant or highly correlated features to simplify the model and reduce overfitting.
3. Cross-Validation: Split the data into training and validation sets, and evaluate the model's performance on the validation set. This helps identify if the model is overfitting and allows for tuning hyperparameters accordingly.

**Exercise 5:** Assumptions of Linear Regression

**Question: What are the assumptions of Linear Regression?**

**Answer:** Linear Regression makes several assumptions about the data:

1. Linearity: The relationship between the independent and dependent variables is linear.
2. Independence: The observations are independent of each other.
3. Homoscedasticity: The variance of the errors is constant across all levels of the independent variables.
4. Normality: The errors follow a normal distribution.
5. No multicollinearity: The independent variables are not highly correlated with each other.

Violations of these assumptions can affect the accuracy and reliability of the regression model. It's essential to check these assumptions and consider appropriate remedies if any assumptions are violated.