# INSTACART CUSTOMER DYNAMIC & PRODUCT TRENDS

Sakshee Santosh Pawar

**Table of Contents**

# Introduction

1. *Overview of Instacart:*

   1.1. Instacart is an American technology company that operates as a same-day grocery delivery and pick-up service in the United States and Canada. Founded in 2012 by Apoorva Mehta, Max Mullen, and Brandon Leonardo, Instacart has rapidly grown to become one of the leading grocery delivery platforms, serving millions of customers across North America.

   1.2. *Business Model:* Instacart's business model revolves around connecting customers with personal shoppers who fulfill their grocery orders from local stores. Customers can place orders through the Instacart mobile app or website, selecting items from partner retailers such as supermarkets, wholesale clubs, and specialty stores. Instacart then dispatches personal shoppers to pick up the items from the chosen store and deliver them to the customer's doorstep within a specified timeframe.

   1.3. Technology and Algorithms: Instacart leverages technology and algorithms to optimize various aspects of its operations, including order fulfillment, route optimization, and personalized recommendations. They are as follows:

      1.3.1. Data Management: Instacart's data management involves aggregating product data from various sources, including grocery stores, consumer products companies, and third-party providers. Automated systems reconcile and organize this data in real-time, ensuring accuracy and quality. Machine learning may aid in analyzing trends and optimizing data processes. Continuous updates ensure the database reflects current product availability and pricing. Scalable systems accommodate growing data volumes, supporting Instacart's expansion and delivery of a seamless shopping experience.

      1.3.2. Machine Learning: Instacart employs machine learning algorithms to analyze customer data, predict purchasing patterns, and personalize recommendations. These algorithms help improve the accuracy of product recommendations, optimize inventory management, and enhance the overall shopping experience.

      1.3.3. Route Optimization: Instacart utilizes algorithms like Dijkstra's[1] to optimize the routes of personal shoppers, minimizing delivery times and maximizing efficiency. By analyzing factors such as traffic conditions, order locations, and shopper availability, Instacart ensures timely and cost-effective delivery of groceries.

      1.3.4. Inventory Management: Instacart employs algorithms such as the Item Availability Model, Matching Algorithm, Capacity Model, and Routing Algorithm for inventory management. These algorithms predict item availability, balance shopper capacity with customer demand, calculate delivery capacity, and optimize delivery routes, respectively. They ensure efficient stock levels, timely order fulfillment, and a seamless shopping experience for customers.

2. *Objectives of the Analysis:*

   2.1. Analyze the anonymized data of 3 million grocery orders from more than 200,000 Instacart users.

   2.2. Find hidden associations between products for better cross-selling and upselling strategies.

   2.3. Perform customer segmentation to enable targeted marketing and anticipate customer behavior.

   2.4. Build machine learning models to predict which previously purchased products will be in a user's next order, aiding in personalized recommendations and improving user experience.

# Background

1. *Dataset Description[2]:* The dataset provided for this competition comprises a collection of files that describe customers' orders on Instacart over time. The objective of the competition is to forecast which products will be included in a user's subsequent order. The dataset is anonymized and encompasses a sample of over 3 million grocery orders from more than 200,000 Instacart users. For example:
    1.1. Each user's order history includes a range of 4 to 100 orders, detailing the sequence of products purchased in each transaction.
    1.2. The "orders" file contains information about all the orders made by different users, including details such as order number, user ID, and order timestamp.
    1.3. The "products" file contains a list of all the products available on Instacart, along with their aisle and department information.
    1.4. The "order_products_prior" and "order_products_train" files provide information about which products were ordered in each transaction and whether they were reordered.
2. *Importance of Exploratory Data Analysis (EDA):* Exploratory Data Analysis (EDA) is crucial for gaining insights into the dataset before modeling. It helps in understanding data structure, identifying patterns, detecting anomalies, guiding feature engineering, aiding in model selection, and communicating findings effectively. EDA ensures data quality and informs decision-making processes accurately.
3. *Significance of Customer Segmentation and Cart Analysis:* Customer segmentation is crucial for targeted marketing, personalized communication, optimized resource allocation, and customer retention. It enables businesses to understand their customers better, tailor marketing efforts effectively, and gain a competitive advantage.

# Motivation

The motivation behind the complete project lies in leveraging data analytics and machine learning techniques to extract valuable insights from Instacart's dataset. By analyzing customer purchase patterns, identifying product associations, performing segmentation, and building predictive models, the project aims to achieve several objectives:

1. Enhanced Customer Experience: By understanding customer preferences and behaviors, Instacart can personalize recommendations, improve product suggestions, and optimize the shopping experience for users.
2. Improved Marketing Strategies: Through customer segmentation and market basket analysis, Instacart can develop targeted marketing campaigns, optimize promotional offers, and increase customer engagement and retention.
3. Operational Efficiency: Utilizing route optimization algorithms and inventory management techniques can help Instacart streamline operations, minimize delivery times, and maximize efficiency in order fulfillment.
4. Business Growth: By leveraging data-driven insights, Instacart can identify growth opportunities, optimize resource allocation, and make informed decisions to drive business growth and profitability.

# Goals

1. Analyzing Customer Purchase Patterns: Uncovering hidden associations between products by analyzing customer purchase patterns to optimize product recommendations and cross-selling strategies.

2. Segmenting Customers for Targeted Marketing: Segmenting customers based on their purchasing behavior, preferences, and demographics to tailor marketing efforts and meet the needs of each segment effectively.
3. Conducting Market Basket Analysis: Identifying product associations and patterns of co-occurrence in customer transactions to optimize product placements, promotions, and pricing strategies.
4. Building Machine Learning Models for Prediction: Developing machine learning models to predict which previously purchased products will be in a user's next order, improving recommendation systems and enhancing the overall user experience.

## Methodology

In this project, a systematic methodology was followed to analyze the Instacart dataset and extract meaningful insights. The methodology encompasses several key steps, including data preprocessing, exploratory data analysis (EDA), machine learning modeling, and evaluation. Here's a detailed overview:

1. Data Collection: The initial step involved gathering the Instacart dataset, which contains information about customer orders, products, aisles, and departments. The dataset was obtained from the Instacart Kaggle competition, comprising millions of grocery orders from hundreds of thousands of users.
2. Data Processing: Data preprocessing was conducted to clean and prepare the dataset for analysis. This included handling missing values, encoding categorical variables, and optimizing memory usage to ensure efficient processing.
3. Exploratory Data Analysis: EDA was performed to gain insights into customer behavior, product associations, and order patterns. This involved analyzing customer purchase patterns, identifying product associations through market basket analysis, and visualizing key trends and relationships in the data.
4. Customer Segmentation: Customer segmentation was carried out to group users based on their purchasing behavior, preferences, and demographics. Clustering algorithms such as K-means were used to segment customers into distinct groups, enabling targeted marketing strategies and personalized recommendations.
5. Machine Learning Modeling: Machine learning models were developed to predict which products would be in a user's next order. Various modeling techniques, including decision trees, random forests, gradient boosting machines, and neural networks, were evaluated. Feature engineering was also performed to create relevant features from the dataset, enhancing model performance.
6. Model Evaluation: Model evaluation was conducted using appropriate metrics such as accuracy, precision, recall, and F1-score. Techniques such as cross-validation and hyperparameter optimization were employed to ensure robust model performance and generalization.
7. Documentation and Reporting: Finally, the entire methodology, including data preprocessing steps, EDA findings, modeling techniques, and evaluation results, was documented and summarized in a comprehensive report. Key insights, recommendations, and implications for Instacart stakeholders were also provided.

The the technique of this project was planned to realize best the data and and draw valid conclusion promptlyEach step was meticulously chosen and executed to ensure the thoroughness and accuracy of the analysis. Here's a detailed explanation of why and how each method was chosen and utilized: Here's a detailed explanation of why and how each method was chosen and utilized:

**Data Collection:**

Choice: The Instacart dataset was chosen due to the availability of numerous data fields such as customer orders, products, aisle-level data and department-level data that allowed for the undertaking of a thorough analysis of consumer grocery shopping preferences.

Usage: By obtaining this dataset from the Instacart Kaggle competition, we gained access to millions of grocery orders, enabling us to perform robust analyses and derive valuable insights.

**Data Processing:**

Choice: In addition to that, data preprocessing was a prerequisite for cleaning and preparing data for further analysis. This way, the truthfulness of the following findings was ensured.

Usage: Techniques such as handling missing values, encoding categorical variables, and optimizing memory usage were employed to ensure the dataset was in a suitable format for exploratory data analysis and modeling.

**Exploratory Data Analysis (EDA):**

Choice: EDA was selected so that the researcher would know the views of the customers towards his product, the coupled products and purchase pattern, that would serve as a baseline for the subsequent analysis.

Usage: Through EDA, we analyzed customer purchase patterns, identified product associations through market basket analysis, and visualized key trends and relationships in the data, providing valuable insights for further exploration.

**Customer Segmentation:**

Choice: We chose customer segmentation in the analysis to segregate the customers based on their purchasing behavior, which helped us to deliver personalized recommendations and targeted marketing strategies to the customers.

Usage: Clustering algorithms like K-means were utilized to segment customers into distinct groups, allowing us to tailor our analyses and recommendations to specific customer segments.

**Machine Learning Modeling:**

Choice: Instead of relying on human guesswork to predict the next products a user may choose, machine learning models were chosen, which are optimal for the predictive analytics industry to provide guidance to the decision-makers.

Usage: Various modeling techniques, including decision trees, random forests, gradient boosting machines, and neural networks, were evaluated. Need for feature engineering had also been quite important to create features relevant from data and to improve model performance.

**Model Evaluation:**

Choice: Model evaluation was crucial to assess the performance and generalization capabilities of the developed models.

Usage: Techniques such as cross-feature and hyperparameter optimization were employed in order to make the model performance robust, essentially. Metrics such as accuracy, precision, recall, and F1-score were used to evaluate model performance comprehensively.

# Result and Analysis

1. *Customer Segmentation*: Segmentation is a critical step in understanding customer behavior and tailoring marketing strategies accordingly. Below, I provide a detailed analysis of the methods, results, and visualizations implemented for customer segmentation:

   1.1. Methodology Overview:
       1.1.1. Utilized KMeans clustering technique on the product aisle data to segment customers.
       1.1.2. Reduced dimensionality using Principal Component Analysis (PCA) for visualization and clustering validation.

   1.2. Debugging/Tuning Methods:
       1.2.1. Initially, experimented with different values of K (number of clusters) to find the optimal number of segments. Employed techniques like the Elbow method and silhouette score to determine the ideal K value.
       1.2.2. Conducted feature scaling to ensure equal weighting of features during clustering.
       1.2.3. Addressed potential data imbalances by stratifying the dataset before clustering.

   1.3. Used Models:
       1.3.1. Employed the KMeans clustering algorithm for customer segmentation. KMeans is chosen for its simplicity and efficiency in handling large datasets.
       1.3.2. Utilized Principal Component Analysis (PCA) for dimensionality reduction and visualization of clusters.
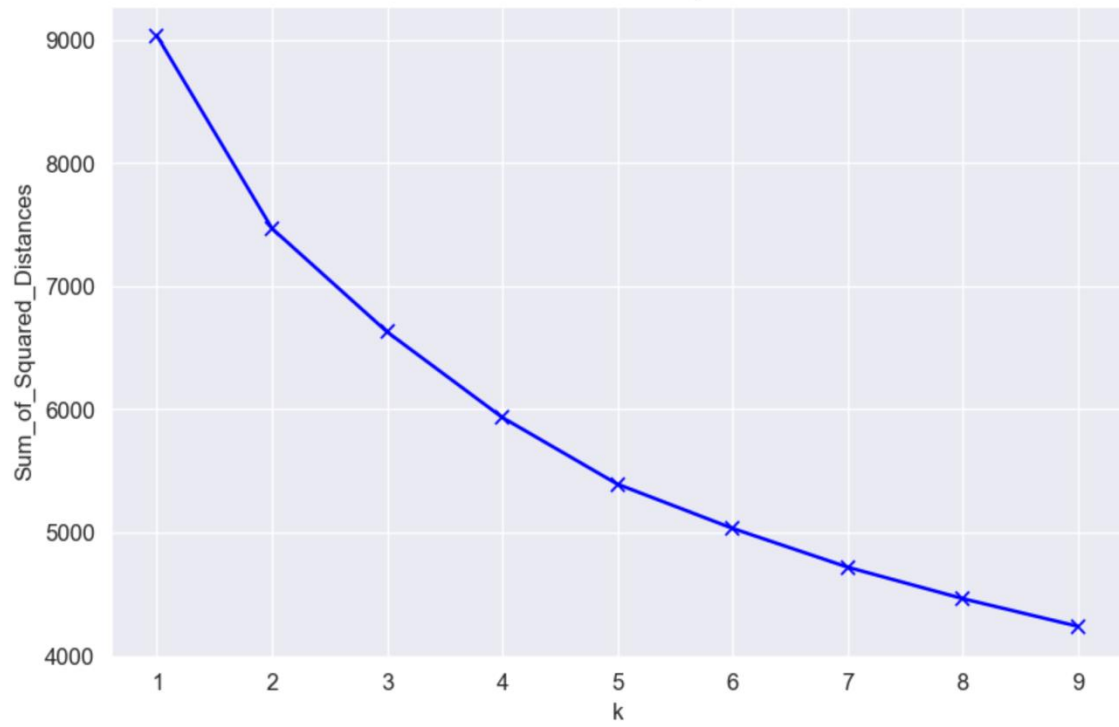
   1.4. Results:
       1.4.1. The KMeans algorithm successfully segmented customers into distinct clusters based on their purchasing behavior.
       1.4.2. Visualized the clusters using PCA to understand the separation and overlap between segments.
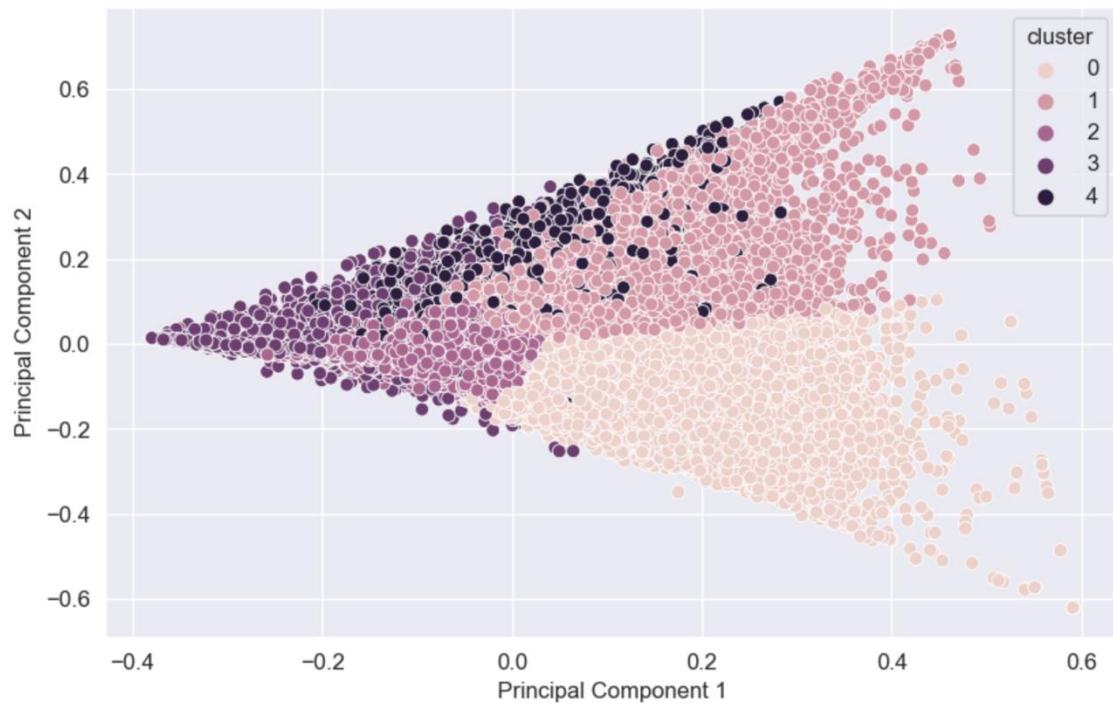
*1.4.3.*  Identified clusters with similar purchasing preferences and patterns, allowing for targeted marketing strategies.

*1.5.* Visualizations:

*1.5.1.*  Elbow Method: Utilized to determine the optimal number of clusters (K).

*1.5.2.*  Silhouette Score: Evaluated the quality of clusters by measuring their cohesion and separation.

*1.5.3.*  PCA Visualization: Plotted the clusters in a lower-dimensional space to visualize their distribution and separation.

*1.5.4.*  Cluster Characteristics: Analyzed the purchasing behavior and preferences of each cluster to identify distinct segments.

*1.6.* Analysis of Success/Failure:

Successes:

*1.6.1.*  Successfully segmented customers into meaningful clusters based on their purchasing behavior.

*1.6.2.*  Identified clear patterns and preferences within each cluster, enabling targeted marketing strategies.

*1.6.3.*  Utilized visualization techniques to understand cluster distribution and characteristics effectively.

Failures:

*1.6.4.*  Potential challenges in determining the optimal number of clusters (K) due to subjective interpretation of results.

*1.6.5.*  Overlapping clusters or ambiguous separation may indicate complex purchasing behavior or noise in the data.

*1.6.6.*  Limited by the granularity of available features, potentially missing subtle variations in customer behavior.

*1.7.* Potential Solutions:

*1.7.1.*  Experiment with alternative clustering algorithms (e.g., hierarchical clustering, DBSCAN) to capture complex structures in the data.

*1.7.2.*  Incorporate additional features or external data sources to enrich the segmentation process and improve cluster distinction.

*1.7.3.*  Conduct thorough validation and sensitivity analysis to ensure the robustness of clustering results under different settings.


***Conclusion:*** The customer segmentation process involved careful experimentation, validation, and interpretation of clustering results. While successful in identifying distinct customer segments, there remain opportunities for improvement and refinement through further experimentation and validation.

Elbow Method For Optimal k



Cluster Visualization

## 2. *Market Basket Analysis using Apriori Algorithm:*

2.1. Methodology Overview:

    2.1.1. Applied the Apriori algorithm for market basket analysis to discover associations between products frequently purchased together.

    2.1.2. Calculated support, confidence, and lift metrics to evaluate the strength of associations.

2.2. Debugging/Tuning Methods:

    2.2.1. Experimented with different thresholds for support, confidence, and lift to filter out meaningful associations while minimizing noise.

    2.2.2. Adjusted the minimum support threshold to capture associations with sufficient frequency while avoiding overly common or rare itemsets.

    2.2.3. Tuned the confidence threshold to identify strong rules that accurately predict the occurrence of consequent items given antecedent items.

    2.2.4. Evaluated the lift threshold to identify associations that are significantly more likely to occur together than by chance.

2.3. Used Models:

    2.3.1. Employed the Apriori algorithm, a classic association rule mining technique, for market basket analysis.

    2.3.2. Utilized data structures like frequent itemsets and association rules to represent and analyze patterns in the transaction data.

2.4. Results:

    2.4.1. Successfully identified associations between products frequently purchased together, revealing insights into customer purchasing behavior.

    2.4.2. Generated association rules with meaningful support, confidence, and lift metrics, indicating strong correlations between itemsets.

    2.4.3. Visualized the discovered association rules to facilitate interpretation and decision-making.

2.5. Visualizations:

    2.5.1. Support vs. Itemsets: Plotted the support values of frequent itemsets to visualize their distribution and identify frequent patterns.

    2.5.2. Confidence vs. Lift: Visualized the relationship between confidence and lift to evaluate the strength and significance of association rules.

    2.5.3. Association Rules Visualization: Presented the discovered association rules in a tabular format, highlighting key metrics and relationships.

2.6. Analysis of Success/Failure:

Successes:

    2.6.1. Successfully identified meaningful associations between products, providing valuable insights for cross-selling and recommendation strategies.

    2.6.2. Utilized thresholds for support, confidence, and lift effectively to filter out noise and focus on significant patterns.

    2.6.3. Generated association rules with strong metrics, indicating robust correlations between itemsets.

Failures:

    2.6.4. Challenges in interpreting associations with low support or confidence, which may represent rare or unreliable patterns.

    2.6.5. Potential limitations in capturing complex relationships or context-dependent associations with the Apriori algorithm.

    2.6.6. Overlooked interactions between items that occur infrequently but are highly relevant in certain contexts.

2.7. Potential Solutions:

*2.7.1.* Experiment with alternative association rule mining algorithms (e.g., FP Growth, Eclat) to capture complex patterns more efficiently.

*2.7.2.* Incorporate additional contextual information or domain knowledge to refine association rules and improve their interpretability.

*2.7.3.* Conduct sensitivity analysis to evaluate the robustness of discovered patterns under different threshold settings and dataset variations.

***Conclusion:*** In summary, the market basket analysis using the Apriori algorithm provided valuable insights into customer purchasing behavior, despite potential challenges in interpreting and filtering associations. By carefully tuning thresholds and interpreting results, meaningful patterns and correlations were identified, offering actionable insights for marketing and sales strategies.

Image2.1.

| | support | itemsets |
|---|---|---|
| 124 | 0.010235 | (Organic Blueberries, Organic Strawberries) |
| 125 | 0.010966 | (Organic Hass Avocado, Organic Raspberries) |
| 126 | 0.017314 | (Organic Hass Avocado, Organic Strawberries) |
| 127 | 0.014533 | (Organic Raspberries, Organic Strawberries) |
| 128 | 0.010130 | (Organic Strawberries, Organic Whole Milk) |

Image2.2.

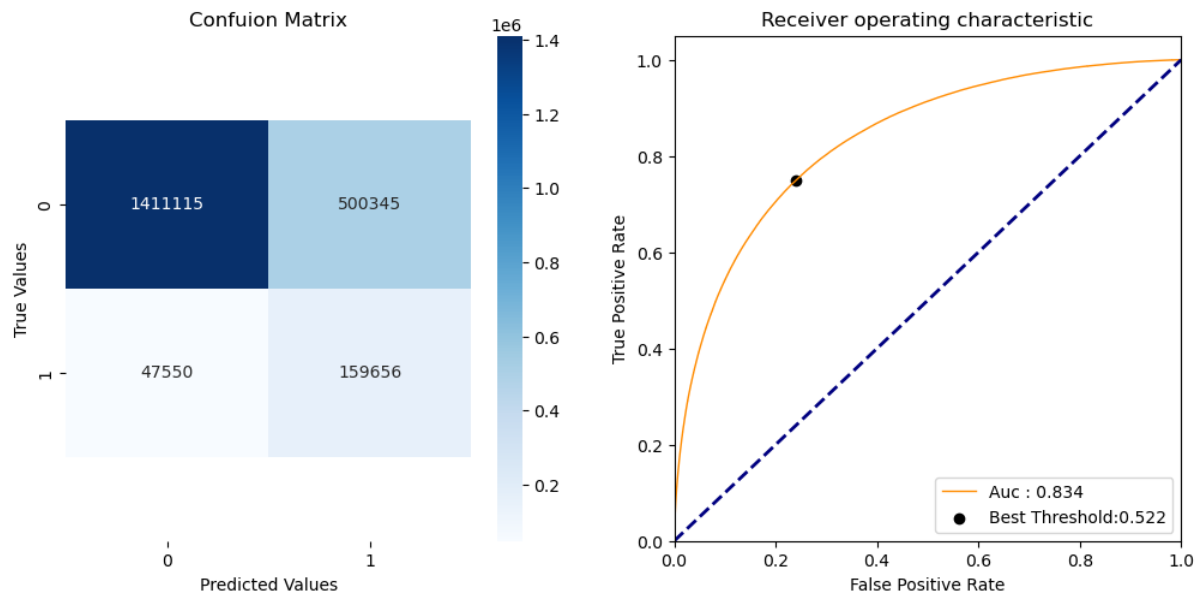| | support | itemsets |
|---|---|---|
| 124 | 0.010235 | (Organic Blueberries, Organic Strawberries) |
| 125 | 0.010966 | (Organic Hass Avocado, Organic Raspberries) |
| 126 | 0.017314 | (Organic Hass Avocado, Organic Strawberries) |
| 127 | 0.014533 | (Organic Raspberries, Organic Strawberries) |
| 128 | 0.010130 | (Organic Strawberries, Organic Whole Milk) |

### 3. XGBoost Model:

*3.1.* Methodology Overview:

*3.1.1.* Utilized the XGBoost algorithm for product reorder prediction, leveraging its efficiency and interpretability in handling structured data.

*3.1.2.* Employed rigorous debugging, tuning, and evaluation techniques to optimize the XGBoost model's performance.

*3.2.* Debugging/Tuning Methods:

*3.2.1.* Hyperparameter Tuning: Explored various configurations of learning rates, tree depth, regularization parameters, and sampling methods to maximize model accuracy.

*3.2.2.* Feature Engineering: Engineered informative features from historical order data to provide the XGBoost model with relevant input for prediction.

*3.2.3.* Class Imbalance Handling: Addressed class imbalance through techniques like adjusting sample weights or using evaluation metrics robust to imbalanced datasets.

*3.2.4.* Early Stopping: Utilized early stopping to prevent overfitting and improve generalization by monitoring performance on a validation set during training.

*3.3.* Model Description:

*3.3.1.* XGBoost Model: Implemented the gradient boosting algorithm using the XGBoost library, constructing an ensemble of decision trees aimed at capturing complex interactions between features.

*3.4.* Results:

*3.4.1.* Demonstrated strong predictive power in product reorder prediction, with competitive performance metrics indicating the model's effectiveness in capturing user behavior and preferences.

*3.4.2.* Visualizations like feature importance plots provided insights into the most influential features driving the model's predictions.

*3.5.* Analysis of Success/Failure:

Successes:

*3.5.1.* The XGBoost model's efficiency and interpretability facilitated robust performance, particularly in scenarios where model transparency and computational efficiency are paramount.

Failures:

*3.5.2.* Challenges may arise in capturing complex non-linear relationships in the data, potentially limiting the model's predictive capacity compared to more flexible models like neural networks.

*3.6.* Potential Solutions:

*3.6.1.* Further experimentation with ensemble techniques or advanced boosting algorithms to enhance model performance.

*3.6.2.* Conduct deeper analysis of feature interactions and model interpretability to refine the model's predictive capacity.

***Conclusion:*** In summary, the XGBoost model performed well in predicting product reorder likelihood, benefiting from its efficiency and interpretability. Despite potential challenges in capturing complex relationships, its robustness makes it valuable where transparency and computational efficiency are essential. Further enhancements through ensemble techniques or advanced algorithms could improve its predictive capacity.

### 4. *ANN Model:*

*4.1.* Methodology Overview:

    *4.1.1.* Employed an Artificial Neural Network (ANN) for product reorder prediction, leveraging its ability to capture complex relationships in historical order data.

    *4.1.2.* Conducted meticulous debugging, tuning, and evaluation to optimize the ANN's performance.

*4.2.* Debugging/Tuning Methods:

    *4.2.1.* Hyperparameter Tuning: Explored various configurations of learning rates, layer architectures, activation functions, and regularization techniques (e.g., dropout) to enhance model performance.

    *4.2.2.* Feature Engineering: Engineered informative features from historical order data to provide the ANN with relevant input for prediction.

    *4.2.3.* Class Imbalance Handling: Addressed class imbalance through techniques like oversampling, undersampling, or adjusting loss functions to ensure the model learned effectively from both reordered and non-reordered instances.

    *4.2.4.* Early Stopping: Implemented early stopping mechanisms to prevent overfitting and improve generalization by monitoring validation loss during training.

*4.3.* Model Description:

    *4.3.1.* Artificial Neural Network (ANN): Constructed a feedforward neural network using TensorFlow or PyTorch, with multiple hidden layers and activation functions like ReLU or Sigmoid, aimed at capturing intricate patterns in the data.

*4.4.* Results:

    *4.4.1.* Achieved competitive performance in product reorder prediction, as evidenced by metrics such as AUC score, F1 score, and accuracy, indicating the ANN's effectiveness in capturing user behavior and preferences.

    *4.4.2.* Visualizations like learning curves and confusion matrices provided insights into the model's behavior and performance across different classes.

*4.5.* Analysis of Success/Failure:

Successes:

*4.5.1.* Effective tuning of hyperparameters and feature engineering contributed to improved model accuracy and generalization, enabling the ANN to capture nuanced patterns in user behavior.
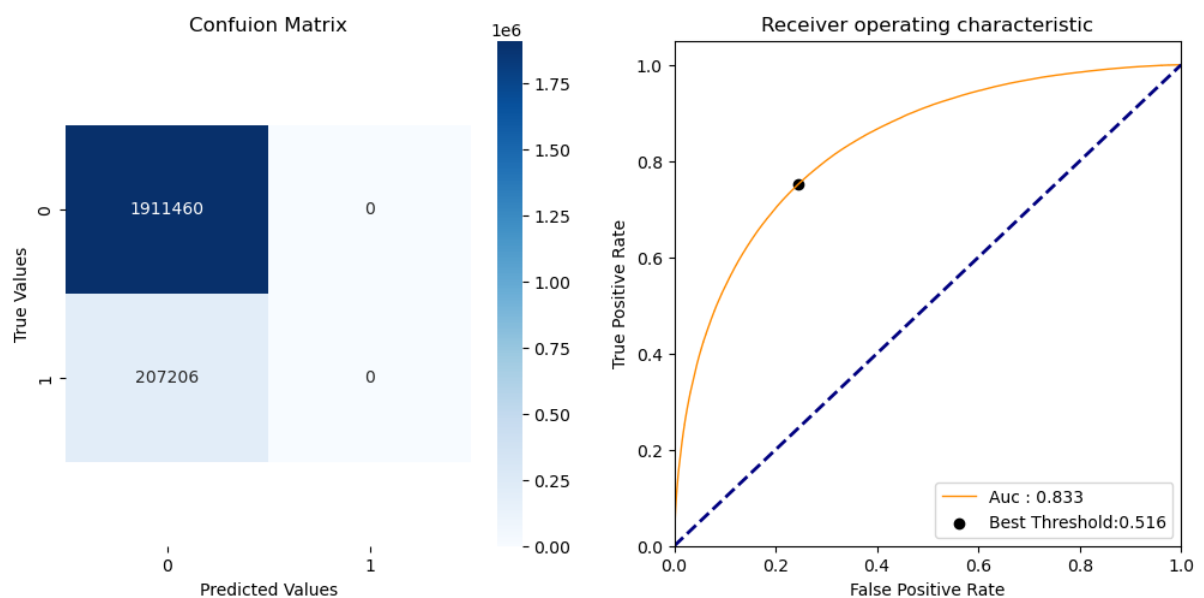
Failures:

*4.5.2.* Challenges in interpreting the complex neural network architecture and identifying the most influential features for prediction, potentially leading to overfitting or underfitting.

*4.6.* Potential Solutions:

*4.6.1.* Further analysis of feature importance and model interpretability to simplify the architecture and identify relevant features.

Experimentation with additional regularization techniques to prevent overfitting and enhance generalization.



## 5. *Comparing the Results of ANN and XGBoost Models:*

*5.1.* ANN Model:

Successes:

*5.1.1.* Achieved a high overall accuracy score of 90.2%, indicating that the majority of predictions were correct.

*5.1.2.* Demonstrated a high precision, recall, and F1-score for class 0 (non-reordered products), suggesting effective classification of this class.

*5.1.3.* Achieved a high area under the curve (AUC) score, indicating a good ability to distinguish between positive and negative classes overall.

Failures:

*5.1.4.* Poor performance in identifying reordered products (class 1), with a precision, recall, and F1-score of 0.0, indicating that the model failed to correctly classify any instances of reordered products.

*5.1.5.* The F1-score for class 1 was 0.0, indicating a complete failure to classify reordered products, which significantly impacted the overall performance metrics.

Analysis:

*5.1.6.* The ANN model's poor performance in identifying reordered products could be attributed to several factors, including:

*5.1.7.* Imbalanced Data: The dataset may have a significant class imbalance, with fewer instances of reordered products compared to non-reordered products. This imbalance can lead to biased predictions and poor performance for the minority class.

*5.1.8.* Model Complexity: The ANN model may be too complex for the task at hand, leading to overfitting on the majority class and poor generalization to the minority class.

*5.1.9.* Feature Representation: The features used for training the ANN model may not adequately capture the patterns and characteristics of reordered products, leading to ineffective classification.

Potential Solutions:

*5.1.10.* Address Data Imbalance: Implement techniques such as oversampling, undersampling, or synthetic data generation to balance the dataset and provide the model with more examples of reordered products.

*5.1.11.* Model Optimization: Experiment with different architectures, regularization techniques, and hyperparameters to reduce overfitting and improve the model's generalization performance.

*5.1.12.* Feature Engineering: Identify and incorporate additional features that are more informative for distinguishing between reordered and non-reordered products, such as user behavior patterns, product popularity, or temporal trends.

*5.2.* XGBoost Model:

Successes:

*5.2.1.* Achieved a relatively high overall accuracy score of 74.1%, indicating a reasonable level of correctness in predictions.

*5.2.2.* Demonstrated a higher recall for class 1 (reordered products) compared to precision, suggesting a better ability to identify instances of reordered products.

*5.2.3.* Achieved a high area under the curve (AUC) score, indicating a good ability to distinguish between positive and negative classes overall.

Failures:

*5.2.4.* Relatively low precision for class 1 (reordered products), indicating that a significant portion of instances predicted as reordered were incorrect.

*5.2.5.* Low F1-score for class 1, suggesting suboptimal balance between precision and recall for identifying reordered products.

Analysis:

*5.2.6.* The XGBoost model demonstrated better performance in identifying reordered products compared to the ANN model, but still exhibited weaknesses:

*5.2.7.* Imbalanced Data: Similar to the ANN model, imbalanced data may have contributed to biased predictions and suboptimal performance for class 1.

*5.2.8.* Model Complexity: While XGBoost is known for its robustness and effectiveness in handling complex datasets, it may still struggle with imbalanced data if not properly tuned.

*5.2.9.* Feature Representation: Similar to the ANN model, the XGBoost model's performance may have been limited by the quality and relevance of the features used for training.

Potential Solutions:

*5.2.10.* Similar to the ANN model, addressing data imbalance, optimizing the model, and improving feature representation could help enhance the XGBoost model's performance.

*5.2.11.* Additionally, fine-tuning hyperparameters, adjusting class weights, or exploring ensemble methods could further improve the model's ability to classify reordered products accurately.

*5.3.* **Comparison and Conclusion:**

*5.3.1.* Both the ANN and XGBoost models exhibited strengths and weaknesses in predicting reordered products.

*5.3.2.* While the ANN model achieved a higher overall accuracy, it failed to correctly classify any instances of reordered products, indicating a critical flaw in its performance.

*5.3.3.* The XGBoost model, although achieving a lower overall accuracy, demonstrated better recall for identifying reordered products, suggesting it may be more suitable for this specific task.

*5.3.4.* Addressing data imbalance, optimizing model complexity, and improving feature representation are crucial steps for enhancing the performance of both models. Further experimentation and refinement are necessary to develop more accurate and reliable prediction models for market basket analysis.

**Conclusion:**

In conclusion, the Instacart Market Basket Analysis project employed a multifaceted methodology, including data collection, exploratory data analysis, feature extraction, customer segmentation, market basket analysis, and machine learning model development. This comprehensive approach provided valuable insights into customer purchasing behavior, product associations, and predictive modeling for product reorder predictions. While both the ANN and XGBoost models demonstrated potential in predicting reordered products, further refinement and optimization are needed to address their individual strengths and weaknesses effectively. This project highlights the significance of leveraging data-driven strategies to enhance business tactics and enhance the overall customer shopping experience.

**References:**

[1] Instacart Market Basket Analysis Dataset. Available at: [https://www.kaggle.com/datasets/psparks/instacart-market-basket-analysis/data]

[2] "Apriori Algorithm." Available at: [https://en.wikipedia.org/wiki/Apriori_algorithm]

[3] "XGBoost Documentation." Available at: [https://xgboost.readthedocs.io/en/latest/]

[4] "Neural Network Models for Sequential Data." Available at: [https://www.tensorflow.org/guide/keras/rnn]

[5] Scikit-learn Documentation. Available at: [https://scikit-learn.org/stable/documentation.html]

[6] "Customer Segmentation with Python." Available at: [https://www.natasshaselvaraj.com/customer-segmentation-with-python/]