# Predicting Consumer Spending Patterns in Retail

Shadin Chatila
DSC 680

White Paper

## Business Problem

In today's competitive retail environment, understanding consumer spending patterns is very critical for businesses to succeed. Companies that can effectively analyze shopping behavior gain a competitive oversight by optimizing marketing strategies, enhancing customer retention, and increasing their profitability as a business.

This project aims to explore key factors influencing customer purchase behavior, focusing on variables such as demographics, product categories, and transaction amounts. By leveraging this data, the goal is to predict spending trends, segment customers, and provide actionable insights to inform targeted marketing campaigns. These insights can empower businesses to personalize customer experiences, allocate resources effectively, and increase sales revenue.

## Background/History

Retail analytics has progressed from your traditional descriptive techniques to advanced predictive analytics techniques. In the past, businesses usually relied on demographic and simple segmentation approaches. However, with access to larger datasets and machine learning tools, businesses can now uncover deeper patterns in consumer behavior. This project builds on these current advancements to create insights that can empower businesses to design targeted marketing campaigns and optimize their operations.
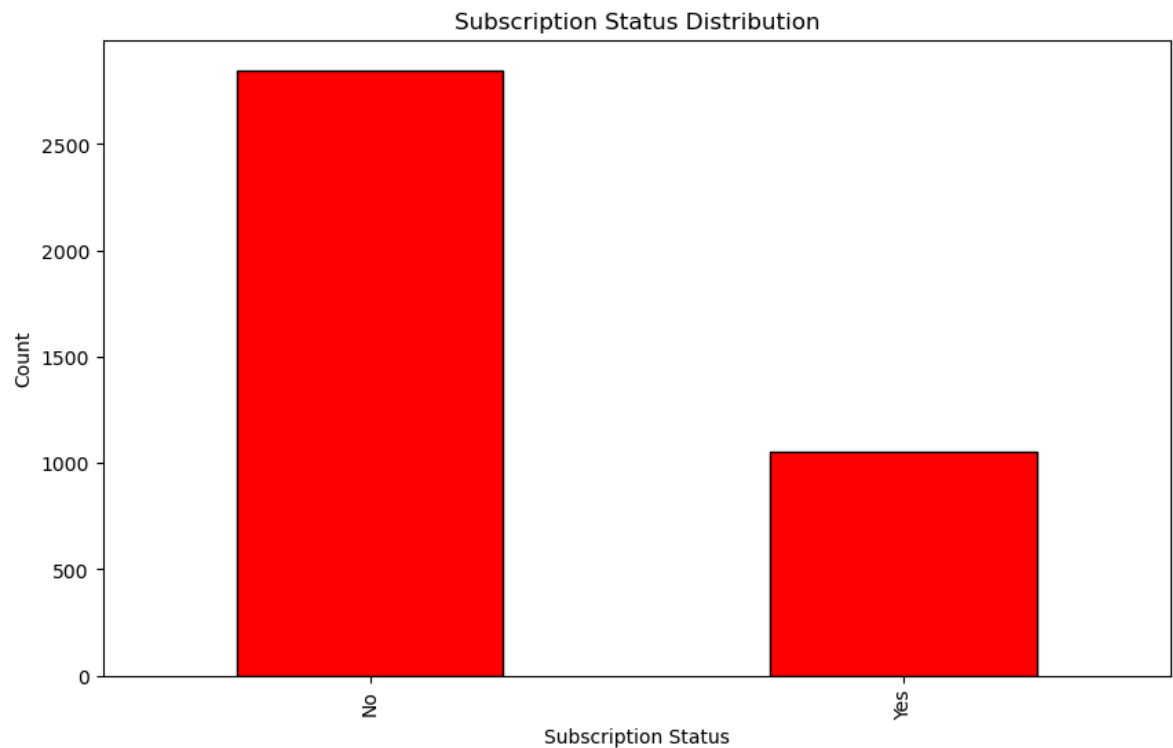
## Data Explanation

The dataset, shopping_behavior.csv, was chosen because it contained a large amount of information on customer shopping activities. Some key variables included are customer demographics, purchase frequency, transaction amounts, product categories, and timestamps of purchases. These features provide a strong foundation for building predictive models and understanding consumer behavior.
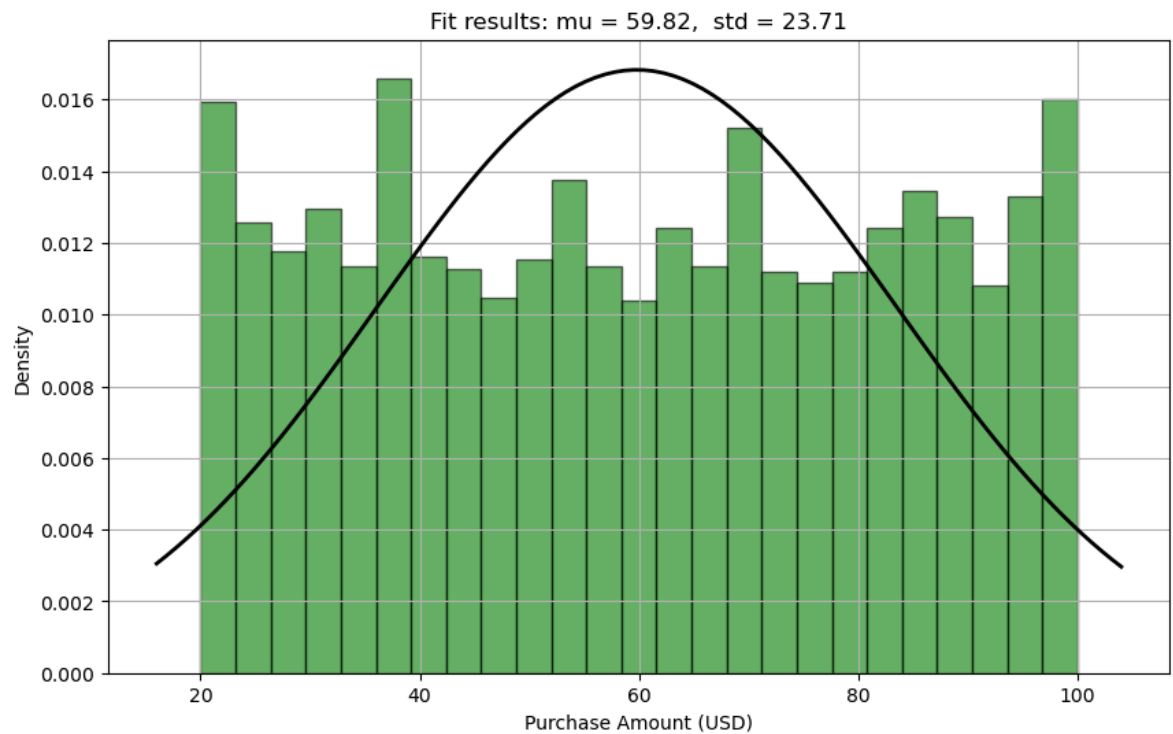
Data preparation involved several critical steps. Missing values were addressed using mean imputation for numerical fields and mode imputation for categorical ones. Duplicate records were removed to ensure data integrity, and categorical variables, such as gender and payment method, were encoded to convert them into a machine compatible format. Numerical features, like the transaction amounts, were standardized to enhance model performance. Finally, the data was split into training and testing subsets using an 80-20 split to ensure robust model evaluation.
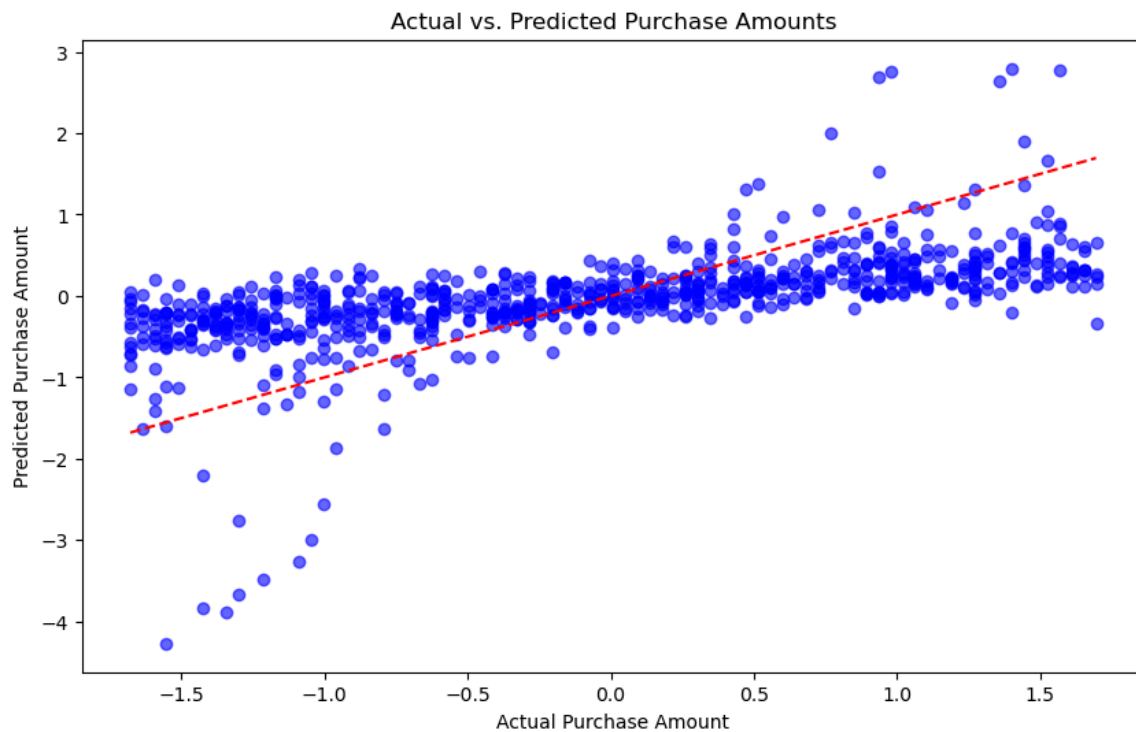
## Methods

Exploratory Data Analysis (EDA) was conducted to understand the distribution and relationships between key variables. A bar plot created from the data shows that 70% of the customers are unsubscribed, which proves to be a significant opportunity for targeted marketing.



Purchase amount distribution was analyzed using a histogram, which showed a mean of $59.82 and a standard deviation of $23.71.

The primary predictive model used was linear regression. The model was preferred as it is commonly used with transactional data, as well as its simplicity and interpretability in predicting continuous outcomes. Key features identified during EDA, such as age, purchase frequency, and subscription status, were used in the model. Model performance was assessed using Mean Squared Error (MSE) and R-squared (R²).



Actual vs. Predicted Purchase Amounts

## Analysis

The initial analysis showed that the older customers generally spend more per transaction, while younger customers tend to make purchases a lot more frequently. Subscription status proved to be a very critical factor, with subscribers displaying more consistent spending habits as one would expect with a subscription model. Additionally, product categories such as electronics significantly influenced higher transaction values. These results align with anticipated trends and highlight the predictive value of demographic and behavioral data.

## Conclusion

The analysis covered highlights the importance of demographic and shopping behavior data in predicting consumer spending within the retail sector. While the linear regression model provides valuable insights, its limitations suggest the need for more advanced modeling techniques to be used. These enhancements can better account for the non-linear relationships involved and help focus on some of the external factors influencing spending behavior.

## Assumptions

This analysis assumes that the dataset accurately represents a typical retail customer base. Additionally, it is assumed that the demographic and behavioral factors, such as age and subscription status, are stable over time and consistently influence consumer spending.

## Limitations

One limitation is the absence of external factors like seasonal trends or economic conditions, which could impact the retail environment and affect consumer spending. Additionally, linear regression may oversimplify relationships, which reduces its ability to capture outlier behaviors or complex interactions.

## Challenges

The main challenges that were faced are related to missing data and maintaining model accuracy. Addressing these data quality issues is critical however time consuming. Also, the simplicity of the linear regression made it a bit challenging to completely capture the complexities of consumer behavior.

## Future Uses/Additional Applications

Incorporating other datasets, such as macroeconomic indicators or holiday sales trends, to enhance the analysis. Clustering techniques could also be applied to segment customers more effectively, enabling more personalized marketing strategies.

## Recommendations

Based on the analysis, businesses should design targeted marketing campaigns for high-value customers, such as offering premium products to the older demographics. Younger frequent shoppers could be incentivized through loyalty programs. Enriching datasets with external variables, such as promotional events or competitor pricing, could also help improve the predictions.

## Implementation Plan

Finalizing the model with cross-validation techniques and then getting it integrated into existing CRM tools. Regular monitoring is conducted to ensure accuracy and will ensure the model remains relevant as consumer behaviors evolve.

## Ethical Assessment

This project adheres to strict data privacy standards, ensuring compliance with regulations and ethical practices. Bias in demographic-driven insights will be mitigated in order to avoid any

unfair targeting. Transparency regarding the model's assumptions and limitations will ensure stakeholders are well-informed about the analysis.

## Appendix

**Figure 1**: Distribution of subscription statuses (referenced in **Methods**)

**Figure 2**: Histogram of transaction amount distributions (referenced in **Methods**)

**Figure 3**: Scatterplot comparing actual and predicted transaction amounts (referenced in **Methods**).

## Questions

**What are the key predictors of spending behavior in your dataset?**
The key predictors that were identified in the dataset included customer demographics (age and gender), subscription status, and product categories.

**How does subscription status impact consumer spending patterns?**
Subscription status significantly impacts spending patterns. Subscribers display more consistent and predictable spending habits, whereas the non-subscribers show higher variability. The data reveals that 70% of customers are unsubscribed, highlighting a potential opportunity for targeted marketing campaigns.

**Are there any notable trends in spending by product category?**
Yes, the electronics and high-value product categories contribute to higher transaction amounts. These categories are consistently associated usually with increased spending, indicating their importance in influencing overall consumer behavior.

**What limitations did the linear regression model present in your analysis?**
The linear regression model presented limitations in capturing the complex, non-linear relationships within the data. Additionally, it also struggled to handle outliers effectively, which reduced its overall predictive accuracy for certain consumer segments.

**How could external factors improve the model?**
Including external factors, such as seasonal trends, macroeconomic conditions, and promotional events, could significantly improve the model. For instance, capturing holiday sales data could refine predictions during peak retail periods.

**Were there any biases observed in the data, and how were they addressed?**
The data exhibited potential biases related to age and subscription status, as these variables heavily influenced spending patterns. To address this, standardization was applied to numerical variables, and careful preprocessing ensured that no single variable disproportionately skewed the model results.

**What future applications do you envision for this model?**
Future applications include integrating the model into CRM systems for personalized marketing, using clustering techniques for customer segmentation, and applying the insights to optimize inventory management and promotional planning.

**How would this analysis scale for larger or global datasets?**
Scaling the analysis for larger or global datasets would require enhanced computational resources and potentially more sophisticated algorithms, such as gradient boosting or deep learning.

**What additional data sources could improve your model's predictions?**
Data sources such as weather data, real-time promotional event tracking, and competitor pricing could enhance the model's ability to account for external factors influencing consumer behavior.

**How were ethical considerations accounted for in the project?**
Ethical considerations were central to the project. Data privacy was maintained by excluding personally identifiable information and ensuring compliance with data protection standards. Biases in the data were carefully mitigated to prevent unfair targeting or discrimination, and all findings were transparently communicated with a clear acknowledgment of the model's limitations and assumptions.

## References

1. Zeesolver. (n.d.). *Consumer behavior and shopping habits dataset*. Kaggle. Retrieved December 1, 2024, from https://www.kaggle.com/datasets/zeesolver/consumer-behavior-and-shopping-habits-dataset