UNIVERSITY OF DOHA FOR SCIENCE AND TECHNOLOGY

# HIGH SPENDER PREDICTION USING MACHINE LEARNING AND SHAP EXPLAINABILITY

MACHINE LEARNING - 3201 PRESENTATION BY
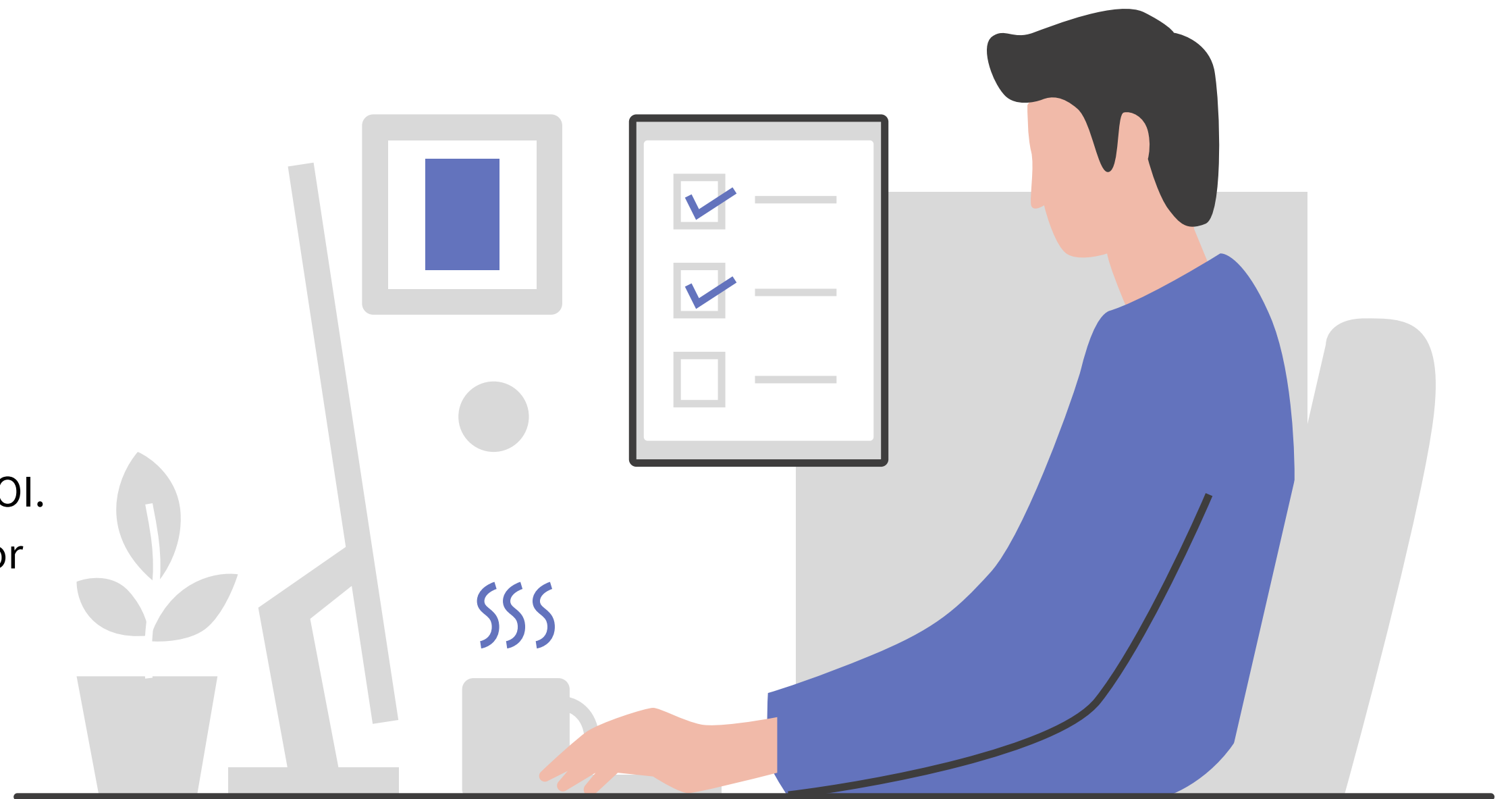
SYED ABDULLAH - 60104641
SHAIMAH MOHAMMED -

# OVERVIEW OF THE PROJECT

## PROJECT OBJECTIVE

In this project, my goal was to analyze customer behavior and predict high-spending customers using customer and purchase data. By identifying these high-value customers, businesses can optimize their marketing efforts, personalize campaigns, and ultimately improve their ROI. Understanding customer spending behavior is crucial for targeted marketing, customer retention, and resource allocation.

LETS GET STARTED

# DATASET DESCRIPTION & BACKGROUND CONTEXT

## BUSINESS PROBLEM

The goal is to predict high-spending customers based on their purchase behavior and demographics. Identifying high spenders allows businesses to target them with personalized offers, optimize marketing efforts, and improve customer retention strategies.

## DATASET OVERVIEW

The dataset is sourced from Kaggle and consists of three files: customer_data.csv, product_data.csv, and purchase_data.csv. It contains a total of 40 columns and approximately 15,000 rows. The dataset includes customer demographics, product details, and transaction history. The main objective of the analysis is to predict high-spending customers based on their purchase behavior. This will help optimize marketing strategies and improve customer retention efforts.

## HOW THE PROBLEM IS ADDRESSED

By merging these datasets, I built a comprehensive dataset to analyze customer behavior. After cleaning and engineering relevant features, I used machine learning models to predict high spenders, helping businesses focus resources on high-value customers.

# WHY I CHOSE THIS DATASET

I chose this dataset because it provides a comprehensive view of customer behavior, product details, and purchasing patterns, which directly aligns with my objective to predict high-spending customers.

## RELEVANCE
The dataset includes customer demographics (e.g., age, gender, region) and purchase history, essential for identifying spending patterns.

## COMPREHENSIVE
It contains both numerical (e.g., total spend) and categorical (e.g., product category, gender) features, providing a solid foundation for analysis.

## DATA QUALITY
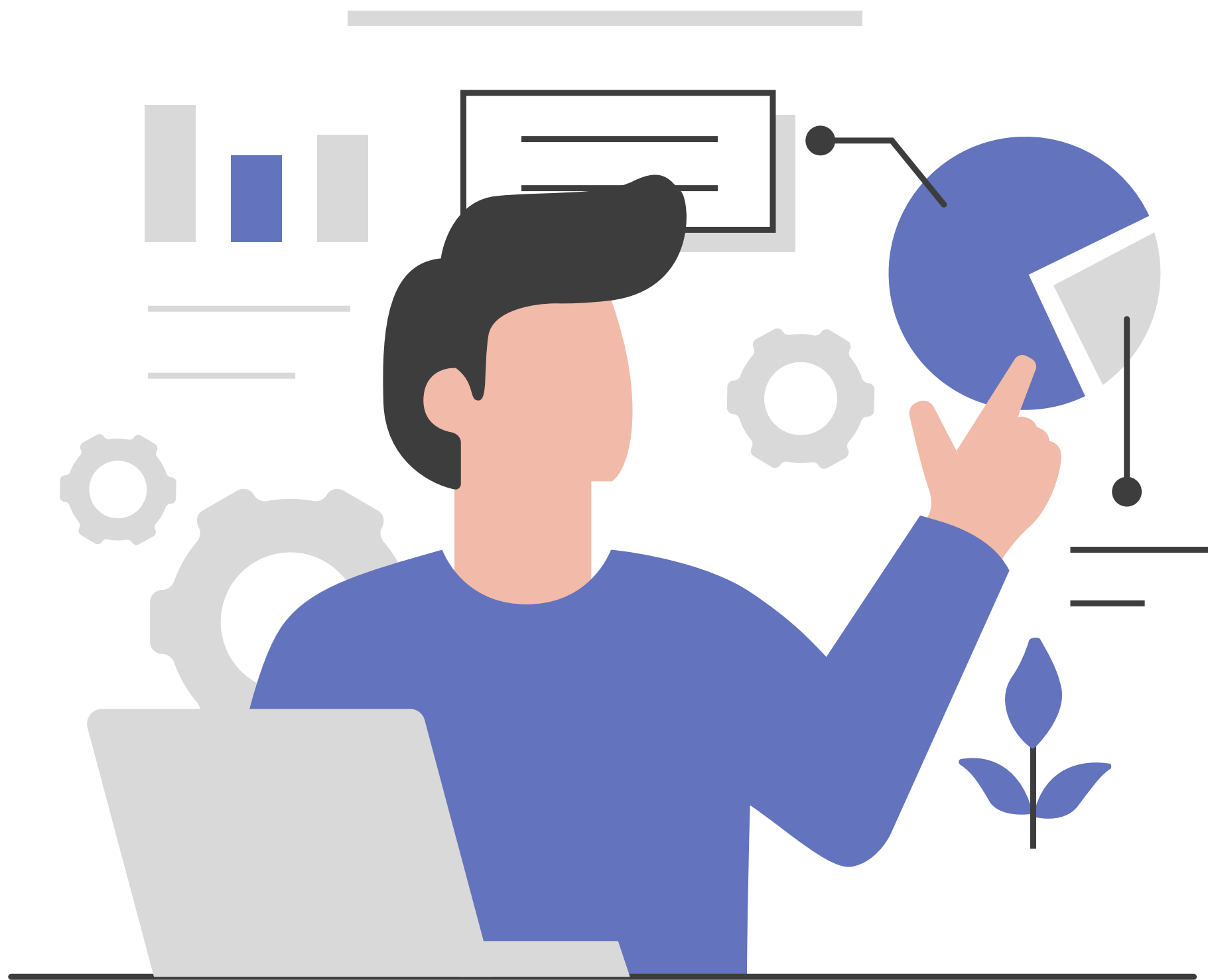The data was well-structured and cleaned, with minor preprocessing required, making it suitable for model building

# VISUAL DATASET SNAPSHOT

## SAMPLE OF DATASET (CUSTOMER, PRODUCT, AND PURCHASE DATA) WITH A FEW ROWS TO PROVIDE CLARITY ON THE STRUCTURE.

| | purchase_id | customer_id | product_id | purchase_date | quantity | total_amount | first_name | last_name | gender | date_of_birth | ... | signup_date | address | city | state |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 42 | 2018-04-15 14:08:01 | 3 | 37.642074 | Robert | Smith | Female | 1994-06-14 21:40:27 | ... | 2016-10-16 17:23:25 | 8465 Main St | San Antonio | CA |
| 1 | 2 | 1 | 138 | 2022-07-10 23:33:47 | 4 | 70.247106 | Robert | Smith | Female | 1994-06-14 21:40:27 | ... | 2016-10-16 17:23:25 | 8465 Main St | San Antonio | CA |
| 2 | 3 | 1 | 403 | 2021-12-31 03:53:33 | 3 | 89.168896 | Robert | Smith | Female | 1994-06-14 21:40:27 | ... | 2016-10-16 17:23:25 | 8465 Main St | San Antonio | CA |
| 3 | 4 | 1 | 193 | 2017-01-14 01:25:11 | 2 | 59.705059 | Robert | Smith | Female | 1994-06-14 21:40:27 | ... | 2016-10-16 17:23:25 | 8465 Main St | San Antonio | CA |
| 4 | 5 | 1 | 26 | 2018-04-06 11:01:06 | 3 | 101.778864 | Robert | Smith | Female | 1994-06-14 21:40:27 | ... | 2016-10-16 17:23:25 | 8465 Main St | San Antonio | CA |

5 rows × 22 columns

# APPROACH:

### DATA PREPROCESSING
I cleaned and transformed the raw data, handling missing values, encoding categorical features, and performing necessary feature engineering

### EXPLORATORY DATA ANALYSIS (EDA):
I explored the data visually to identify trends, correlations, and patterns in customer behavior, which helped me select features for the model.

### MODEL BUILDING
I used Logistic Regression and Shap to automatically train and tune machine learning models.

# DATA CLEANING & PREPROCESSING

Before starting the analysis, I performed essential data cleaning and preprocessing steps to ensure the data was ready for modeling:

### Missing Values

I handled missing values by either imputing (filling in) or removing rows/columns with too many missing values

### Duplicate Records

I identified and removed any duplicate entries in the dataset to avoid bias.

### Data Transformation

I created new features like signup years (calculated from signup_date) and total spending per customer to improve the model's predictive power

### Feature Encoding

Categorical features like gender and product category were encoded into numerical values for compatibility with the machine learning model.

# VISUAL DATA CLEANING SNAPSHOT

## DATA CLEANING: A SCREENSHOT OF THE CODE USED FOR HANDLING MISSING VALUES OR ENCODING FEATURES.

```python
print("\nMissing values in Product Data:")
print(product_df.isnull().sum())

print("\nMissing values in Purchase Data:")
print(purchase_df.isnull().sum())

# Drop duplicates
customer_df.drop_duplicates(inplace=True)
product_df.drop_duplicates(inplace=True)
purchase_df.drop_duplicates(inplace=True)
```

```
Missing values in Customer Data:
customer_id       0
first_name        0
last_name         0
gender            0
date_of_birth     0
email             0
phone_number      0
signup_date       0
address           0
city              0
```

# EDA OVERVIEW

EXPLORATORY DATA ANALYSIS (EDA) WAS AN ESSENTIAL PART OF THIS PROJECT, AS IT HELPED ME UNDERSTAND THE DATASET AND UNCOVER MEANINGFUL PATTERNS BEFORE BUILDING THE PREDICTIVE MODEL.

## THE MAIN OBJECTIVES OF THE EDA WERE:

- To understand the distribution of key variables.
- To identify any patterns or correlations between features.
- To visualize how different factors (like age, gender, and region) impact spending behavior.

# EXPLORATORY DATA ANALYSIS

In this section, we perform descriptive statistics, correlations, and visualizations to explore trends and patterns. These insights will guide our modeling and business recommendations.

```
[196]:   # General statistical summary
         full_df.describe(include='all')
```

[196]:

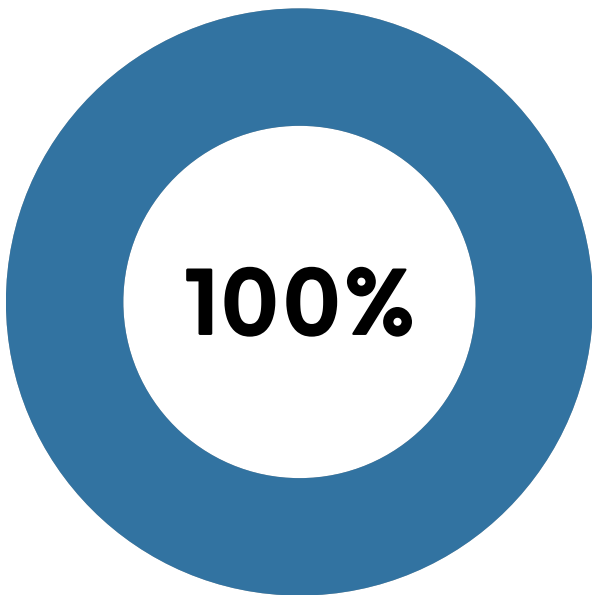| | purchase_id | customer_id | product_id | purchase_date | quantity | total_amount | first_name | last_name | gender | date_of_birth | ... | signup_date | address |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 10308.000000 | 10308.000000 | 10308.000000 | 10308 | 10308.000000 | 10308.000000 | 10308 | 10308 | 10308 | 10308 | ... | 10308 | 10308 |
| unique | NaN | NaN | NaN | 10308 | NaN | NaN | 10 | 10 | 2 | 1000 | ... | 1000 | 990 |
| top | NaN | NaN | NaN | 2018-04-15 14:08:01 | NaN | NaN | Alex | Smith | Female | 1973-02-01 13:43:23 | ... | 2016-06-18 01:41:15 | 7346 Main St |
| freq | NaN | NaN | NaN | 1 | NaN | NaN | 1240 | 1158 | 5496 | 20 | ... | 20 | 28 |
| mean | 5154.500000 | 504.540648 | 251.363795 | NaN | 3.030656 | 77.423841 | NaN | NaN | NaN | NaN | ... | NaN | NaN |
| std | 2975.807621 | 292.026758 | 143.690280 | NaN | 1.412852 | 58.719304 | NaN | NaN | NaN | NaN | ... | NaN | NaN |
| min | 1.000000 | 1.000000 | 1.000000 | NaN | 1.000000 | 1.526648 | NaN | NaN | NaN | NaN | ... | NaN | NaN |
| 25% | 2577.750000 | 245.750000 | 127.000000 | NaN | 2.000000 | 30.143436 | NaN | NaN | NaN | NaN | ... | NaN | NaN |
| 50% | 5154.500000 | 510.000000 | 253.000000 | NaN | 3.000000 | 62.499946 | NaN | NaN | NaN | NaN | ... | NaN | NaN |
| 75% | 7731.250000 | 758.000000 | 375.250000 | NaN | 4.000000 | 113.776378 | NaN | NaN | NaN | NaN | ... | NaN | NaN |
| max | 10308.000000 | 1000.000000 | 500.000000 | NaN | 5.000000 | 249.963513 | NaN | NaN | NaN | NaN | ... | NaN | NaN |

11 rows × 22 columns

# DISTRIBUTION OF PURCHASES BY PRODUCT CATEGORY

In this step of the Exploratory Data Analysis (EDA), I analyzed the distribution of purchases across different product categories to identify which categories generate the highest number of purchases. This helps in understanding customer preferences and pinpointing high-performing product segments.
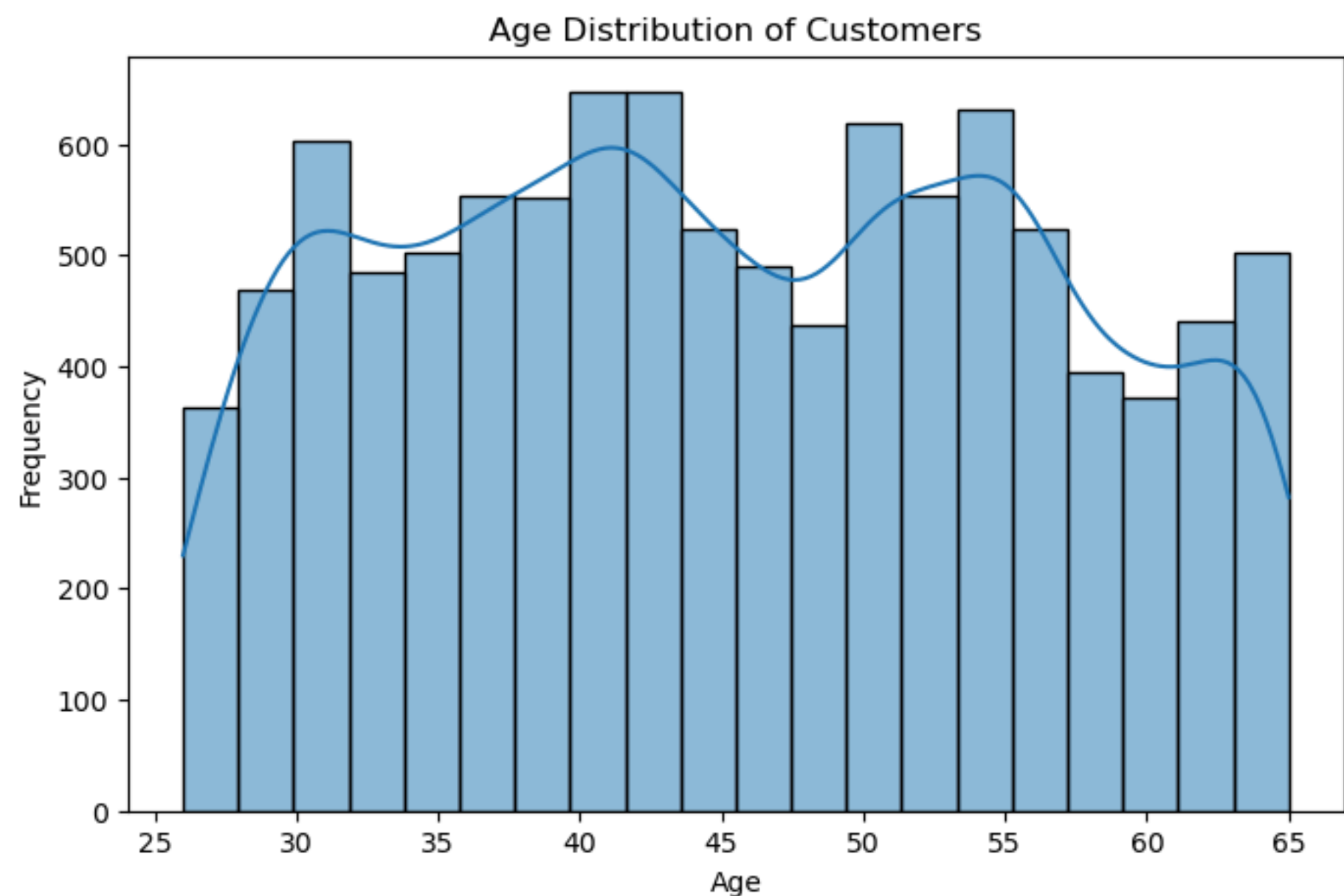


Number of Purchases by Product Category

**ANALYSIS**

I used a horizontal bar chart to visualize the number of purchases in each product category.

**100%**

**INSIGHT**

The analysis revealed that Fruits and Meats are the highest purchased categories, followed by Grains and Dairy. This indicates the potential for businesses to focus on these categories for targeted promotions and marketing efforts.

**100%**

# AGE DISTRIBUTION OF CUSTOMERS



Age Distribution of Customers

**HIGHEST SALES**

**400K**

A histogram with a KDE (Kernel Density Estimate) overlay was used to visualize the distribution of customer ages. The histogram shows the frequency of customers in each age group, while the KDE provides a smooth estimation of the distribution.

**INSIGHTS**

**70%**

The age distribution appears relatively uniform with a slight peak in the 30-45 age range. This suggests that customers in their 30s to 40s form a significant portion of the customer base, and businesses can target this group with age-appropriate products and services.

**CHART SUMMARY**

In this part of the Exploratory Data Analysis (EDA), I examined the age distribution of customers. This helps to understand the demographics of the customer base, which can inform targeted marketing efforts, product recommendations, and customer service

# CUSTOMERS BY TOTAL SPEND

In this analysis, I identified the top 10 customers based on their total spend. This is crucial for understanding high-value customers and can be used to create loyalty programs or targeted marketing campaigns.

## ANALYSIS

The bar chart highlights the total amount spent by the top 10 customers. Each customer is represented by their Customer ID and the corresponding total amount spent.

## INSIGHT

The customers with the highest total spend are likely the most valuable to the business. Businesses should focus on maintaining and nurturing these high-spending customers, as they contribute a significant portion of revenue.

# FEATURE ENGINEERING & TARGET DEFINITION

Feature engineering is the process of transforming raw data into meaningful features to improve machine learning model performance. For this project, we aimed to predict whether a customer is a high spender based on their purchase behavior and profile

## IDENTIFYING RELEVANT FEATURES

Focused on critical features impacting high-spender prediction, such as:
Total Purchase Amount, Product Categories, Age, Signup Year, Gender

## FEATURE CREATION

Aggregating Data: Created new features like Avg Spend and Total Quantity to summarize customer purchase behavior.

## TARGET DEFINITION

HighSpender: Defined as customers with total spend above the 75th percentile.

## DATA PREPROCESSING

- Removed irrelevant columns (purchase_id, purchase_date).
- Encoded categorical features like Category and Gender using One-Hot Encoding.

# SHAP SUMMARY PLOT

## DISTRIBUTION OF IMPACT

### EXPLANATION

This plot shows which features had the most impact on the model's predictions.
 Each dot is one prediction. The color shows the feature value (red = high, blue = low).
 Position on the x-axis shows how much that feature increased or decreased the prediction.

### WHY IT MATTERS

It helps us understand how different features affect predictions, not just how important they are.



**Logistic Regression - SHAP Feature Importance**

| Feature | Value |
|---|---|
| worst texture | +0.96 |
| radius error | +0.8 |
| worst radius | +0.76 |
| worst area | +0.75 |
| worst concave points | +0.72 |
| area error | +0.64 |
| worst perimeter | +0.61 |
| worst smoothness | +0.6 |
| worst concavity | +0.55 |
| worst symmetry | +0.54 |
| mean concave points | +0.54 |
| mean compactness | +0.48 |
| compactness error | +0.43 |
| mean area | +0.42 |
| Sum of 16 other features | +3.62 |

# SHAP FEATURE IMPORTANCE

AVERAGE IMPACT ON MODEL PREDICTION

## EXPLANATION

This bar chart shows the average impact of each feature across all predictions.
The higher the bar, the more important the feature.

## WHY IT MATTERS

It shows which features the model relies on most when predicting high-spending customers.



Random Forest - SHAP Feature Importance

# BIAS ANALYSIS



Total Purchase Amount by Gender (Zoomed In)



Purchase Distribution by Gender

**1** **GENDER COMPARISON**

- We checked if males or females spent more or were predicted differently.
- Results show both genders have similar total spending and almost equal purchase share, with females slightly higher at 53.3%.

**2** **WHY IT'S IMPORTANT**

- To make sure the model is fair and not favoring one gender.
- Our data shows a balanced pattern in total spending and number of purchases between males and females, so no strong gender bias.

Product and Sales Insights Dashboard