

<b>EXPT NO:2</b>	<b>Implementation of data visualization techniques</b>
<b>DATE: 06.01.2026</b>	

### **PRE-LAB QUESTIONS (PROVIDE BRIEF ANSWERS TO THE FOLLOWING QUESTIONS)**

1. **Why is exploratory data analysis critical before model building?**  
It helps understand data structure, detect errors, missing values, patterns, and relationships before choosing or training a model.
2. **How do distributions influence algorithm selection in ML?**  
Some algorithms assume specific distributions (e.g., normality), while skewed or non-linear distributions may require robust or non-parametric models.
3. **What insights can outliers provide in business data?**  
Outliers can indicate anomalies, fraud, rare events, data errors, or high-value opportunities needing special attention.
4. **Why are visual summaries preferred over raw tables?**  
Visuals quickly reveal trends, patterns, and anomalies that are hard to identify from large numerical tables.
5. **How does visualization improve business intelligence?**  
It enables faster decision-making by clearly communicating insights, trends, and performance metrics to stakeholders.

### **IN-LAB EXERCISE:**

#### **OBJECTIVE:**

To explore data distribution and variability using advanced visualization techniques.

#### **SCENARIO:**

A startup analyzes e-commerce transaction data to understand customer spending behavior and detect abnormal purchase patterns.

### **IN-LAB TASKS (Using R Language)**

- Plot histogram of transaction amounts
- Use boxplot to detect outliers
- Create heatmap of monthly sales intensity

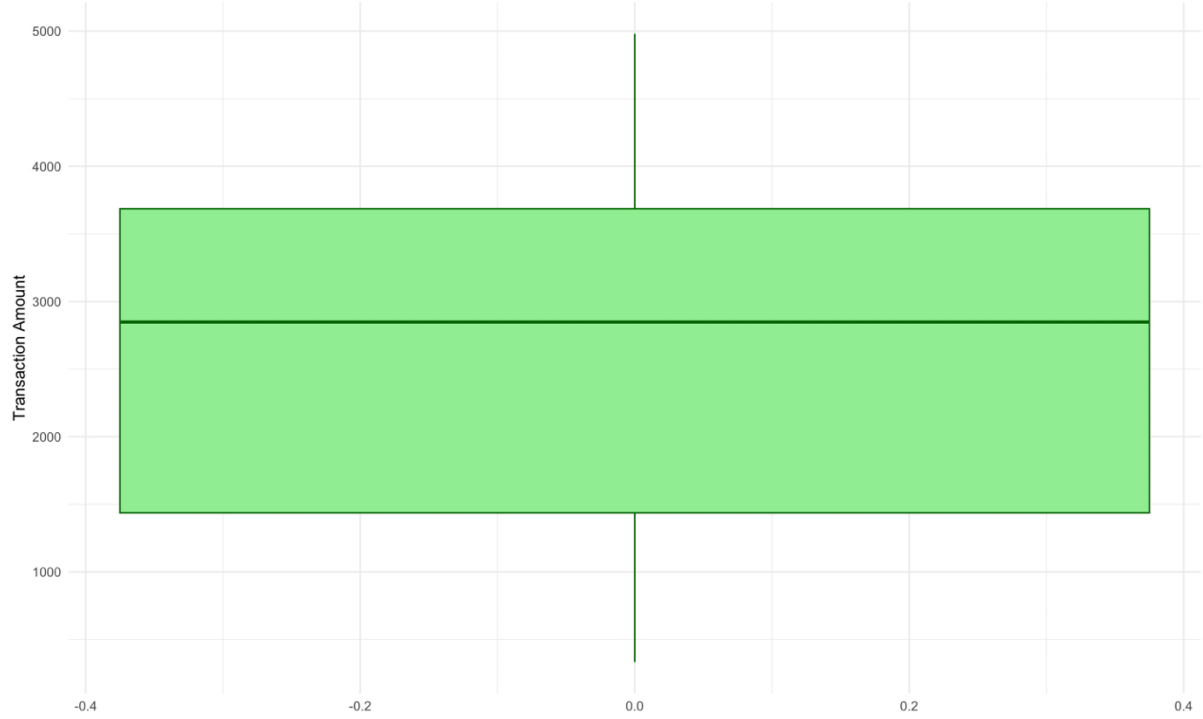
## Code:

Ex-2 Implementation of data visualization techniques

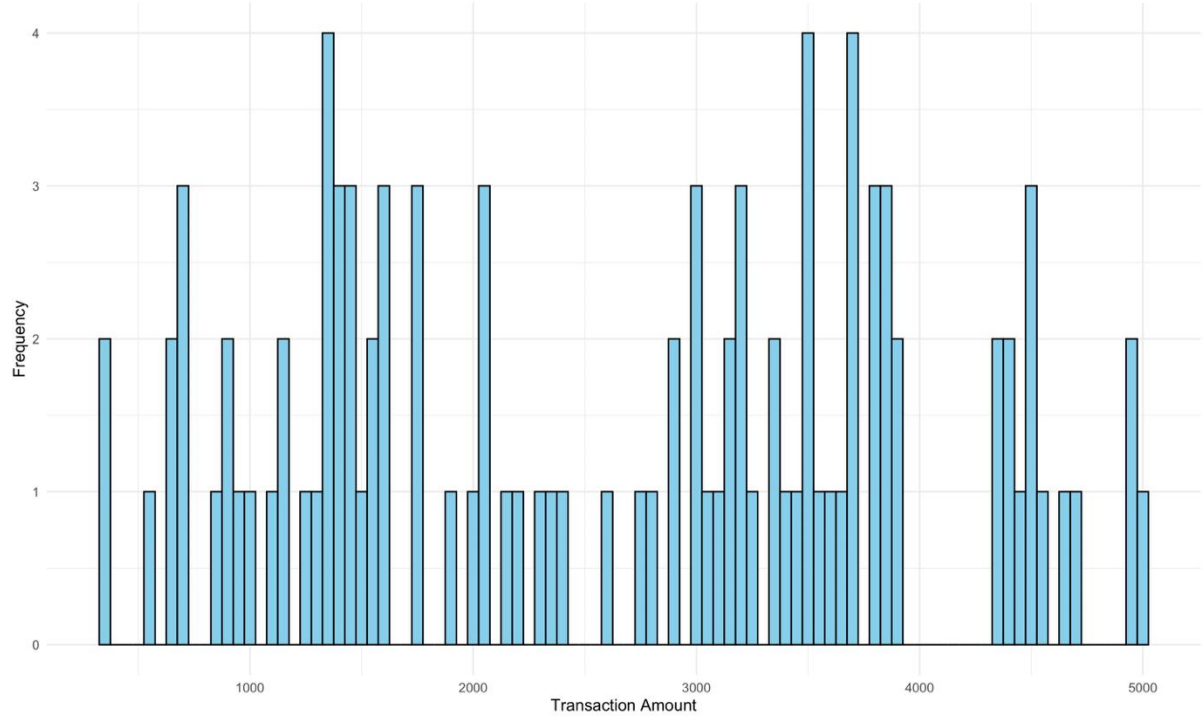
```
1 # Load required libraries
2 library(ggplot2)
3 library(dplyr)
4 # Load your CSV (correct filename)
5 transactions_df <- read.csv("/Users/sugitha/Downloads/2_ecommerce_transactions.csv")
6
7 head(transactions_df)
8
9 # 1 Histogram of Transaction Amounts
10 ggplot(transactions_df, aes(x = Transaction_Amount)) +
11   geom_histogram(binwidth = 50, fill = "skyblue", color = "black") +
12   labs(title = "Histogram of Transaction Amounts",
13        x = "Transaction Amount",
14        y = "Frequency") +
15   theme_minimal()
16
17 # 2 Boxplot to detect outliers
18 ggplot(transactions_df, aes(y = Transaction_Amount)) +
19   geom_boxplot(fill = "lightgreen", color = "darkgreen") +
20   labs(title = "Boxplot of Transaction Amounts",
21        y = "Transaction Amount") +
22   theme_minimal()
23
24 # Extract and print outliers
25 outliers <- boxplot.stats(transactions_df$Transaction_Amount)$out
26 cat("Outliers in Transaction Amounts:\n")
27 print(outliers)
28
29 # 3 Monthly Sales Heatmap
30 # Convert Transaction_Date to Date type
31 transactions_df$Transaction_Date <- as.Date(transactions_df$Transaction_Date, format="%Y-%m-%d")
32
33 # Extract month-year
34 transactions_df$Month <- format(transactions_df$Transaction_Date, "%Y-%m")
35
36 # Aggregate total sales by month
37 monthly_sales <- transactions_df %>%
38   group_by(Month) %>%
39   summarise(Total_Sales = sum(Transaction_Amount), .groups = "drop")
40
41 # Heatmap
42 ggplot(monthly_sales, aes(x = Month, y = 1, fill = Total_Sales)) +
43   geom_tile(color = "white") +
44   scale_fill_gradient(low = "yellow", high = "red") +
45   labs(title = "Monthly Sales Heatmap",
46        x = "Month",
47        y = "",
48        fill = "Total Sales") +
49   theme_minimal() +
50   theme(axis.text.y = element_blank(),
51         axis.ticks.y = element_blank(),
52         axis.title.y = element_blank(),
53         axis.text.x = element_text(angle = 45, hjust = 1)) +
54   coord_fixed(ratio = 0.1)
55
```

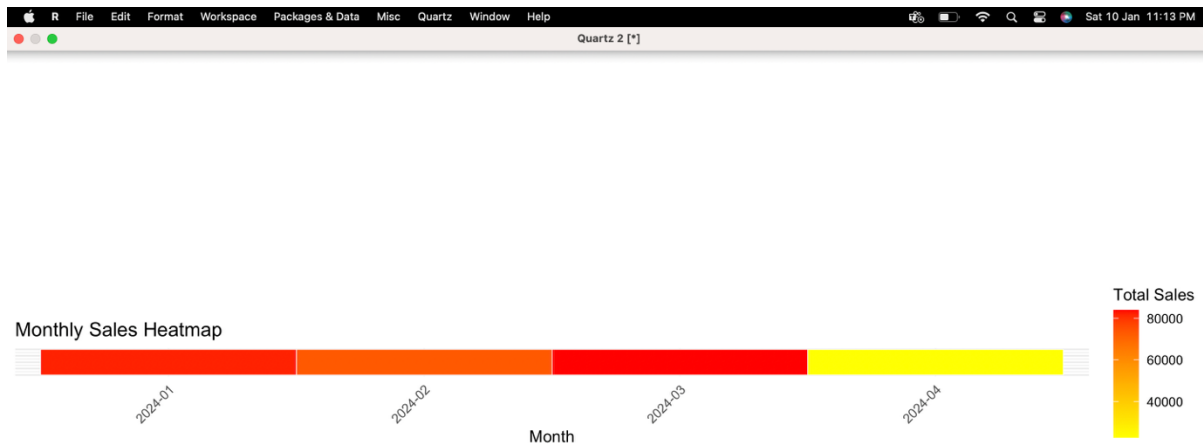
## Output:

Boxplot of Transaction Amounts



Histogram of Transaction Amounts





## POST-LAB QUESTIONS (PROVIDE BRIEF ANSWERS TO THE FOLLOWING QUESTIONS)

### 1. What does a right-skewed distribution indicate about customer behavior?

A right-skewed (positively skewed) distribution of transaction amounts means most customers spend smaller amounts, but a few customers make very large purchases. This suggests that high spenders are rare but can significantly impact revenue.

### 2. How can detected outliers impact business decisions?

Outliers can indicate fraudulent transactions, data entry errors, or unusual high-value purchases. Identifying them helps businesses:

- Investigate anomalies or potential fraud
- Adjust marketing strategies for high-value customers
- Avoid skewed analytics that could misguide decisions

### 3. Which visualization best supports anomaly detection?

Boxplots are particularly effective for anomaly detection because they clearly highlight data points outside the normal range (outliers). Histograms can show general distribution, but boxplots pinpoint unusual values.

### 4. How does EDA improve AI model accuracy?

EDA helps identify:

- Outliers, missing values, and errors

- Feature distributions and correlations
- Skewed or imbalanced data

By cleaning and understanding the data first, AI models learn from high-quality inputs, reducing bias and improving predictive accuracy.

### 5. How can visualization guide feature engineering?

Visualizations reveal patterns, trends, and relationships that may not be obvious from raw tables. For example:

- Scatterplots can suggest interaction features
- Histograms or density plots can guide transformations (e.g., log transform for skewed features)
- Heatmaps reveal correlations, suggesting features to combine or drop

### ASSESSMENT

Description	Max Marks	Marks Awarded
Pre Lab Exercise	5	
In Lab Exercise	10	
Post Lab Exercise	5	
Viva	10	
<b>Total</b>	<b>30</b>	
<b>Faculty Signature</b>		