# DATA  ANALYTICS  MINI  PROJECT

# Title :  Mall Customer Segmentation & Analytics

Name : SUSHANT  GOVIND  SHETTI

Roll No.: 22B-CO-062

## 1. Overview

The primary objective of this project is to develop an interactive dashboard that analyzes customer demographics and spending behavior. By employing unsupervised machine learning (K-Means Clustering), the project aims to identify distinct customer segments (target groups) to aid marketing strategies and business decision-making.

## 2. Dataset Description

Source: Kaggle ("Mall Customer Segmentation Data")

https://www.kaggle.com/datasets/abdallahwagih/mall-customers-segmentation

Size: 200 records

Attributes:

The dataset contains the following 5 variables:

1. CustomerID: Unique ID assigned to the customer.
2. Gender: Gender of the customer (Male/Female).
3. Age: Age of the customer.
4. Annual Income (k$): Annual income of the customer in thousands of dollars.
5. Spending Score (1-100): A score assigned by the mall based on customer behavior and spending nature (1 = Low spender, 100 = High spender).

## 3. Technology Stack & Libraries Used

The project is built using the R Programming Language within the RStudio environment. The following libraries were utilized to achieve specific functionalities:

| Library | Purpose |
| --- | --- |
| **shiny** | The core framework used to build the interactive web application (UI and Server). |
| **bslib** | Used to upgrade the User Interface (UI) to a modern, responsive "Dashboard" style (Bootstrap 5). |
| **bsicons** | Provides modern icons for the Value Boxes (KPI indicators). |
| **dplyr** | Used for data manipulation, cleaning (renaming columns), and filtering. |
| **ggplot2** | The primary engine for creating static visualization charts (Histograms, Boxplots, Scatter plots). |
| **DT** | Renders the raw data as an interactive, searchable HTML table. |
| **corrplot** | Generates the graphical correlation matrix to visualize relationships between variables. |
| **thematic** | Automatically applies the dashboard's CSS theme/colors to the R plots for visual consistency. |

## 4. Analytical Techniques & Objectives

### A. Exploratory Data Analysis (EDA)

Before complex modeling, we perform EDA to understand the data distribution.

- Technique: Histograms and Boxplots.
- Objective:
  - To visualize the frequency distribution of Age, Income, and Spending Score.
  - To analyze gender-based differences (e.g., "Do females generally have a higher spending score than males?").

### B. Correlation Analysis

- Technique: Pearson Correlation Coefficient ($r$).
- Visual: Heatmap Matrix.
- Objective: To quantify the linear relationship between numerical variables.
  - *Example:* We check if there is a strong correlation between Age and Spending Score. (A negative correlation would imply that as customers get older, they spend less).
  - *Key Insight:* In this specific dataset, Correlation Analysis often reveals that Annual Income is *not* linearly correlated with Spending Score, necessitating the use of Clustering.

### C. K-Means Clustering (Unsupervised Machine Learning)

This is the core advanced analytics feature of the project.

- Technique: K-Means Algorithm.
- Objective: To segment customers into distinct groups based on two factors simultaneously: Annual Income and Spending Score.
- How it works:
  - The user selects $k$ (number of clusters) via a slider.
  - The algorithm initiates $k$ centroids.
  - It assigns every customer to the nearest centroid based on Euclidean distance.
  - It iteratively updates the centroids to minimize variance within groups.
- Business Value: This allows the business to identify specific segments, such as:
  - *Careful Spenders:* Low Income, Low Spending.
  - *Standard Customers:* Average Income, Average Spending.
  - *Target Customers (VIPs):* High Income, High Spending.

## 5. User Interface (UI) Features

The application moves beyond standard layouts by utilizing a Modern Dashboard Framework:

1. Value Boxes: Top-level KPIs (Key Performance Indicators) displaying Total Customers, Average Income, and Average Spending Score for immediate insights.
2. Card-Based Layout: Visualizations are contained within distinct "cards" with headers and borders for better organization.
3. Interactive Controls:
   - Dropdowns: To switch between variables dynamically.
   - Sliders: To adjust the $k$ value in K-Means clustering in real-time.
4. Responsive Design: The layout automatically adjusts for different screen sizes (using the flatly Bootstrap theme).

## 6. CODE :

**app.R**

```r
1. library(shiny)
2. library(ggplot2)
3. library(dplyr)
4. library(DT)
5. library(corrplot)
6. library(bslib)      # For modern UI themes and cards
7. library(bsicons)    # For icons in the value boxes
8. library(thematic)   # To make plots match the theme automatically
9.
10.  # --- 1. DATA LOADING & CLEANING ---
11.  df <- tryCatch({
12.    data <- read.csv("Mall_Customers.csv")
13.    colnames(data) <- c("CustomerID", "Gender", "Age",
   "Annual_Income_k", "Spending_Score")
14.    data
15.  }, error = function(e) {
16.    # Fallback data if file is missing (so app doesn't crash during
   testing)
17.    data.frame(
18.      CustomerID = 1:200,
19.      Gender = sample(c("Male", "Female"), 200, replace = TRUE),
20.      Age = sample(18:70, 200, replace = TRUE),
21.      Annual_Income_k = sample(15:137, 200, replace = TRUE),
22.      Spending_Score = sample(1:100, 200, replace = TRUE)
23.    )
24.  })
25.
26.  # Enable thematic for auto-styling plots to match the dashboard
27.  thematic_shiny()
28.
29.  # --- 2. UI SECTION ---
30.  ui <- page_sidebar(
31.
32.    # Set a modern theme (try "flatly", "minty", "journal", or "darkly")
33.    theme = bs_theme(bootswatch = "flatly"),
34.
35.    title = "Customer Intelligence Dashboard",
36.
37.    # The Sidebar
38.    sidebar = sidebar(
39.      title = "Controls",
40.      selectInput("hist_var", "Distribution Variable:",
41.                  choices = c("Age", "Annual_Income_k",
   "Spending_Score")),
42.
43.      sliderInput("clusters", "K-Means Clusters:",
```

```r
44.                   min = 2, max = 6, value = 5),
45.
46.       hr(),
47.       helpText("Use the slider to adjust customer segmentation
   groups."),
48.       p(class = "text-muted", "v2.0 Modern Build")
49.     ),
50.
51.     # The Main Dashboard Area
52.
53.     # ROW 1: Key Performance Indicators (Value Boxes)
54.     layout_columns(
55.       fill = FALSE,
56.       value_box(
57.         title = "Total Customers",
58.         value = nrow(df),
59.         showcase = bs_icon("people-fill"),
60.         theme = "primary"
61.       ),
62.       value_box(
63.         title = "Avg Annual Income",
64.         value = paste0("$", round(mean(df$Annual_Income_k), 1), "k"),
65.         showcase = bs_icon("cash-coin"),
66.         theme = "teal"
67.       ),
68.       value_box(
69.         title = "Avg Spending Score",
70.         value = round(mean(df$Spending_Score), 1),
71.         showcase = bs_icon("graph-up"),
72.         theme = "purple"
73.       )
74.     ),
75.
76.     br(), # Little bit of spacing
77.
78.     # ROW 2: The Main Tabs inside a Card
79.     card(
80.       card_header("Analytics Workspace"),
81.       tabsetPanel(
82.
83.         # Tab 1: Visualization
84.         tabPanel("Distributions",
85.                  layout_columns(
86.                    col_widths = c(6, 6),
87.                    card(full_screen = TRUE, plotOutput("dist_plot")),
88.                    card(full_screen = TRUE, plotOutput("box_plot"))
89.                  )
90.         ),
91.
92.         # Tab 2: Advanced Analysis
93.         tabPanel("Segmentation (AI)",
```

```
 94.                    layout_columns(
 95.                        col_widths = c(5, 7), # 5 parts correlation, 7 parts
    clustering
 96.                        card(
 97.                            card_header("Correlation Matrix"),
 98.                            plotOutput("corr_plot")
 99.                        ),
100.                        card(
101.                            card_header("K-Means Cluster Map"),
102.                            plotOutput("cluster_plot")
103.                        )
104.                    )
105.            ),
106.
107.            # Tab 3: Data Table
108.            tabPanel("Raw Data",
109.                    DTOutput("raw_table")
110.            )
111.        )
112.    )
113. )
114.
115. # --- 3. SERVER SECTION ---
116. server <- function(input, output) {
117.
118.    # 1. Render Data Table
119.    output$raw_table <- renderDT({
120.      datatable(df, options = list(pageLength = 10, scrollX = TRUE),
121.                style = "bootstrap4") # Modern table style
122.    })
123.
124.    # 2. Render Distribution Histogram
125.    output$dist_plot <- renderPlot({
126.      req(input$hist_var)
127.      ggplot(df, aes_string(x = input$hist_var)) +
128.        geom_histogram(bins = 20, fill = "#2C3E50", color = "white") + #
    Matching Flatly theme dark blue
129.        labs(title = paste("Distribution of", input$hist_var), y =
    "Frequency") +
130.        theme_minimal(base_size = 14)
131.    })
132.
133.    # 3. Render Boxplot
134.    output$box_plot <- renderPlot({
135.      req(input$hist_var)
136.      ggplot(df, aes_string(x = "Gender", y = input$hist_var, fill =
    "Gender")) +
137.        geom_boxplot() +
138.        labs(title = paste(input$hist_var, "by Gender")) +
139.        theme_minimal(base_size = 14)
140.    })
```
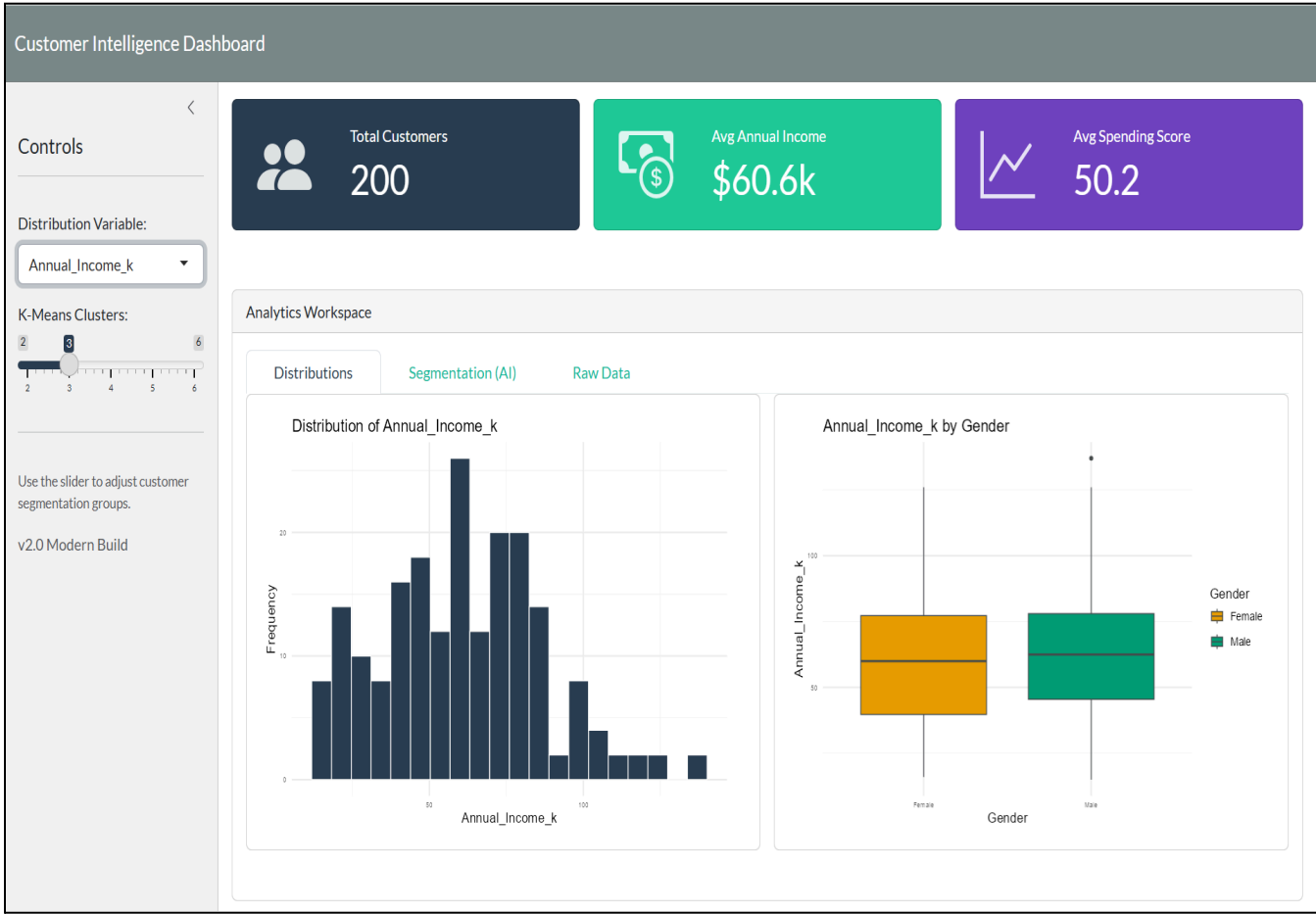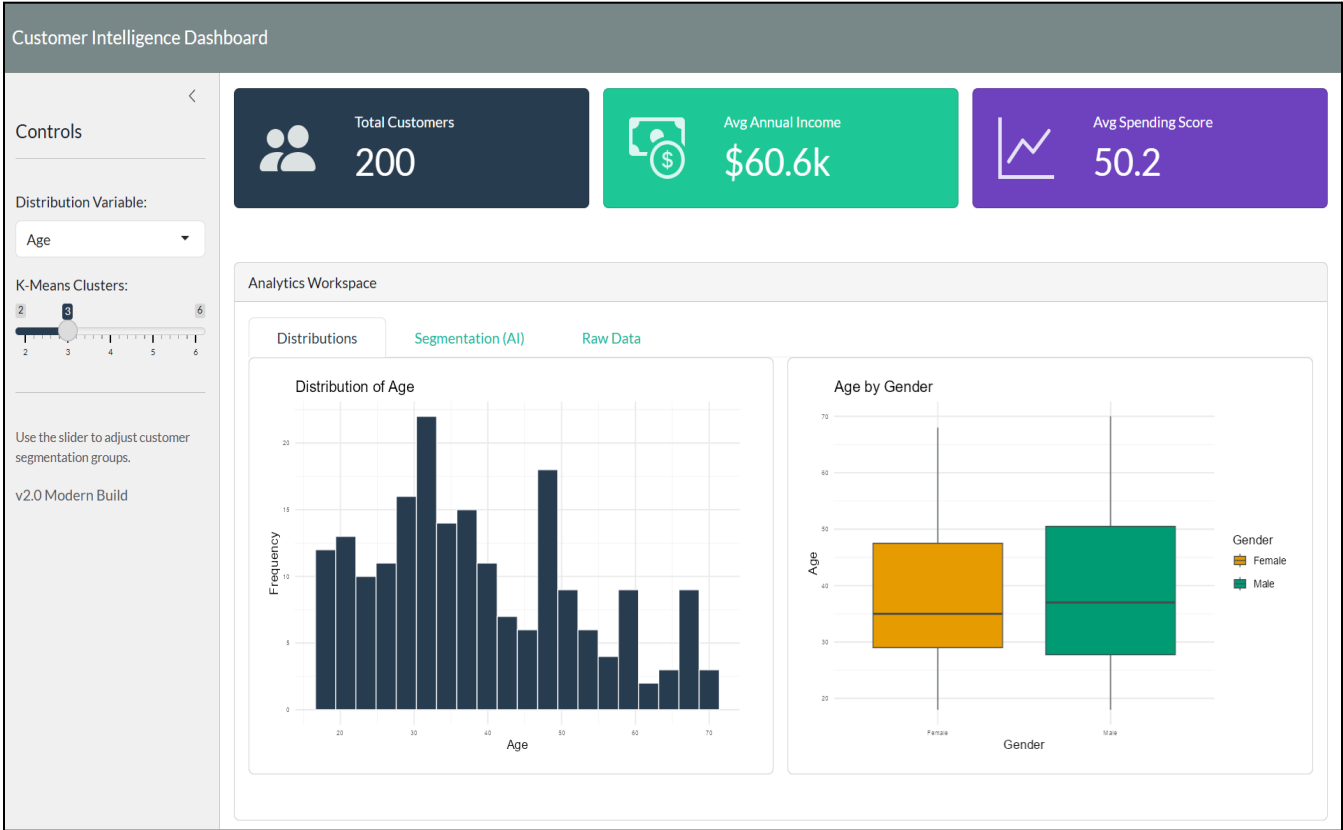
```r
141.
142.    # 4. Render Correlation Plot
143.    output$corr_plot <- renderPlot({
144.      num_data <- df %>% select(Age, Annual_Income_k, Spending_Score)
145.      M <- cor(num_data)
146.      corrplot(M, method = "color", type = "upper",
147.               addCoef.col = "black",
148.               tl.col = "black", tl.srt = 45,
149.               diag = FALSE)
150.    })
151.
152.    # 5. Render Clustering (K-Means)
153.    output$cluster_plot <- renderPlot({
154.      cluster_data <- df %>% select(Annual_Income_k, Spending_Score)
155.      set.seed(123)
156.      kmeans_result <- kmeans(cluster_data, centers = input$clusters)
157.      cluster_data$Cluster <- as.factor(kmeans_result$cluster)
158.
159.      ggplot(cluster_data, aes(x = Annual_Income_k, y = Spending_Score,
  color = Cluster)) +
160.        geom_point(size = 5, alpha = 0.8) +
161.        stat_ellipse(aes(fill = Cluster), geom = "polygon", alpha = 0.2)
  + # Adds fancy circles around clusters
162.        theme_minimal(base_size = 14) +
163.        labs(title = paste("Customer Segments (k =", input$clusters,
  ")"),
164.             x = "Annual Income (k$)",
165.             y = "Spending Score (1-100)") +
166.        scale_color_brewer(palette = "Set1") +
167.        scale_fill_brewer(palette = "Set1")
168.    })
169.  }
170.
171.  # --- 4. RUN APP ---
172.  shinyApp(ui = ui, server = server)
```

**OUTPUT (Screenshots):**

# Customer Intelligence Dashboard

## Controls

### Distribution Variable:
Spending_Score ▼

### K-Means Clusters:
2  3  6
2  3  4  5  6

Use the slider to adjust customer segmentation groups.

v2.0 Modern Build

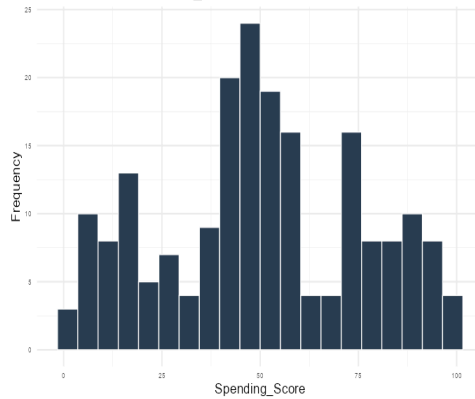| Total Customers | Avg Annual Income | Avg Spending Score |
|---|---|---|
| 200 | $60.6k | 50.2 |

## Analytics Workspace

**Distributions** | Segmentation (AI) | Raw Data

### Distribution of Spending_Score



### Spending_Score by Gender



Gender
■ Female
■ Male

---

# Customer Intelligence Dashboard

## Controls

### Distribution Variable:
Annual_Income_k ▼

### K-Means Clusters:
2  4  6
2  3  4  5  6

Use the slider to adjust customer segmentation groups.

v2.0 Modern Build

| Total Customers | Avg Annual Income | Avg Spending Score |
|---|---|---|
| 200 | $60.6k | 50.2 |

## Analytics Workspace

Show 10 entries                                           Search:

| CustomerID | Gender | Age | Annual_Income_k | Spending_Score |
|---|---|---|---|---|
| 111 | 111 Male | 65 | 63 | 52 |
| 112 | 112 Female | 19 | 63 | 54 |
| 113 | 113 Female | 38 | 64 | 42 |
| 114 | 114 Male | 19 | 64 | 46 |
| 115 | 115 Female | 18 | 65 | 48 |
| 116 | 116 Female | 19 | 65 | 50 |
| 117 | 117 Female | 63 | 65 | 43 |
| 118 | 118 Female | 49 | 65 | 59 |
| 119 | 119 Female | 51 | 67 | 43 |
| 120 | 120 Female | 50 | 67 | 57 |

Showing 111 to 120 of 200 entries

Previous  1  …  11  12  13  …  20  Next

**Customer Intelligence Dashboard**

Controls

Distribution Variable:
Spending_Score ▾

K-Means Clusters:
2    4    6
2  3  4  5  6

Use the slider to adjust customer segmentation groups.

v2.0 Modern Build

Total Customers
200

Avg Annual Income
$60.6k

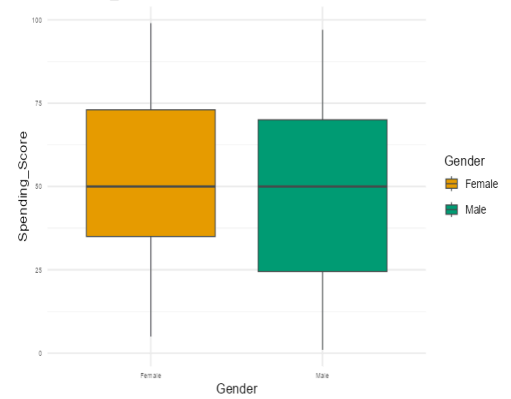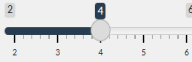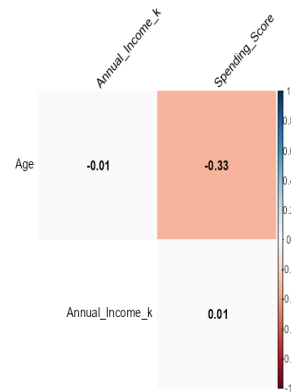Avg Spending Score
50.2

Analytics Workspace

Distributions    Segmentation (AI)    Raw Data

Correlation Matrix

K-Means Cluster Map

Customer Segments (k = 4 )

Cluster
1
2
3
4

# 8. Conclusion

This project successfully transforms raw transactional data into actionable business insights. By integrating statistical analysis with machine learning (Clustering) in an interactive Shiny dashboard, we provide a tool that allows stakeholders to not only visualize past data but also identify high-value customer segments for targeted marketing campaigns.