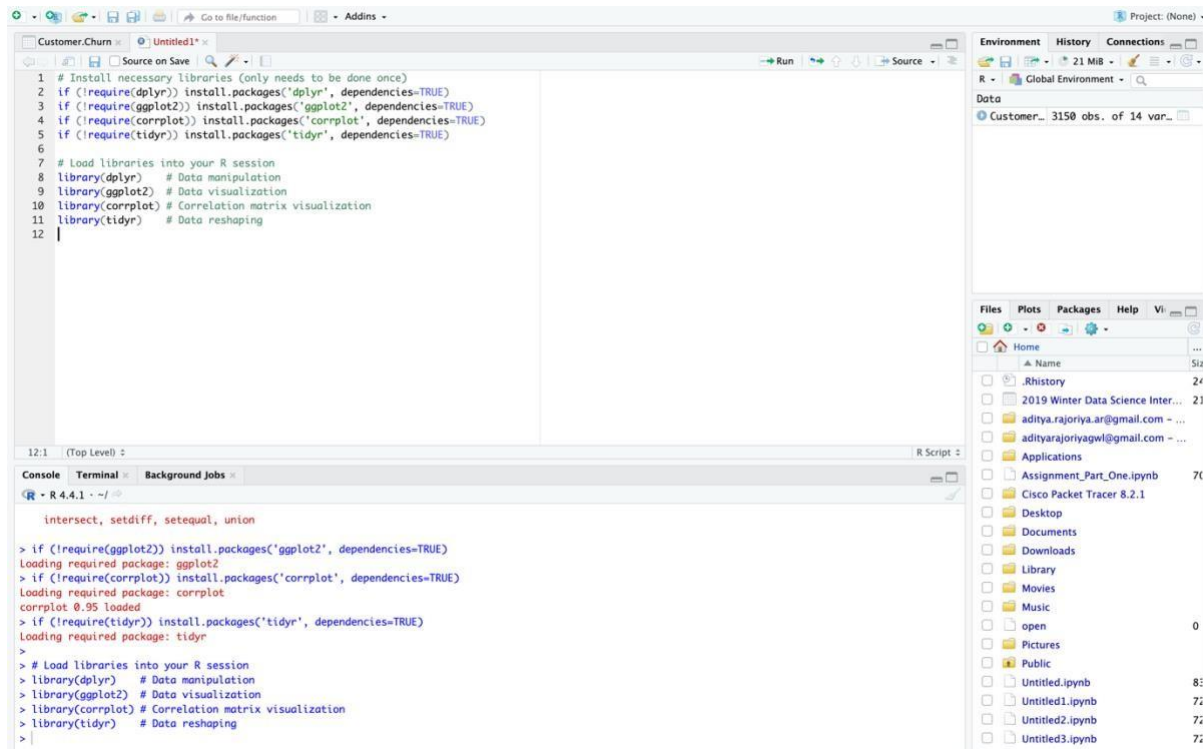Aditya Shroff
101539669

# PROGRAMMING FUNDAMENTALS FOR ANALYTICS

## R_Assignment

**Name: Aditya Shroff (101539669)**
**Professor: Mohsen Selseleh**

Aditya Shroff
101539669

# Step 1: Load the dataset and necessary libraries



# 2. Load the Dataset and Handle Missing Values

Aditya Shroff
 101539669

## 3. Calculate IQR and Remove Outliers



## 4. Visualize Before and After Outlier Removal

Aditya Shroff
101539669

## 5. Group Data and Visualize Average Customer Value by Churn



## 6. Create and Visualize Correlation Matrix
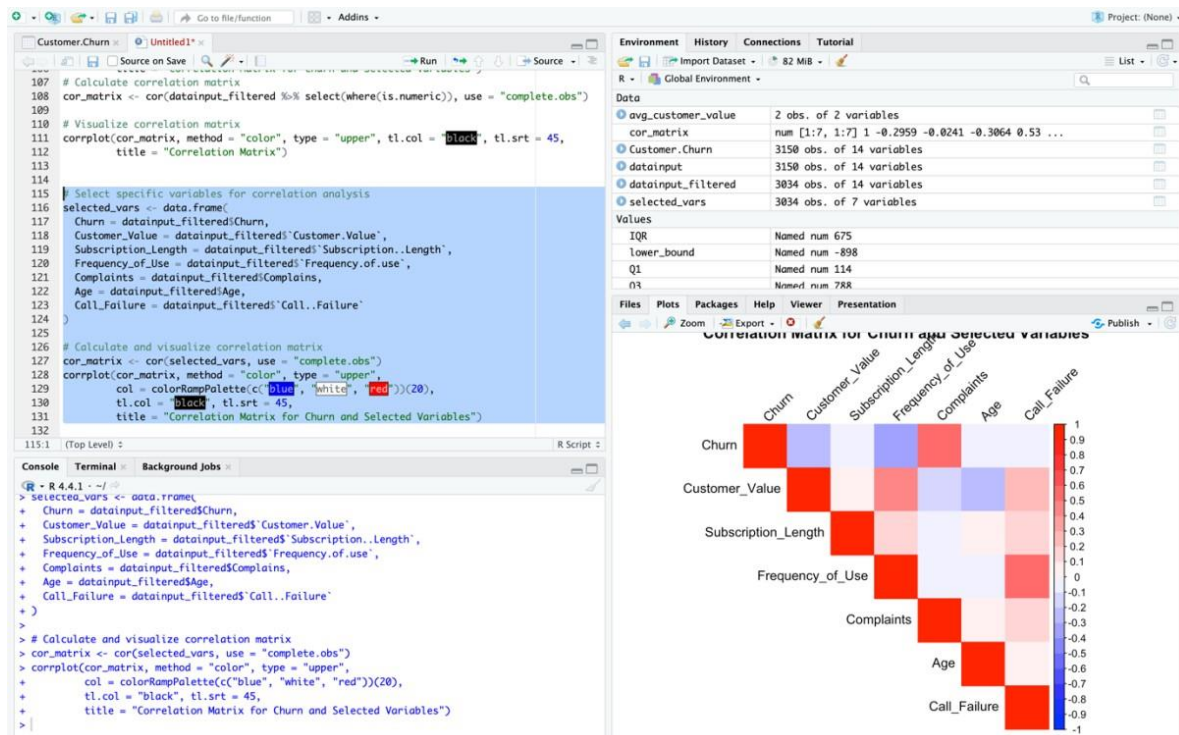
Aditya Shroff
101539669

## 7. Display Column Names



## 8. Select Specific Variables and Create Another Correlation Matrix

## 9. Aggregate Data and Calculate Summary Statistics



## 10. Perform T-Test

## 11. Perform ANOVA



## 12. Build and Summarize Regression Model