Name:Chavan Dhananjay Rajendra
Div:TE3 Roll No:6
Batch:A

# EXPERIMENT-11

EXPT11 : Implement Linear regression using R tool
Theory:  Linear regression is a regression model that uses a straight line to describe the relationship between variables. It finds the line of best fit through given data by searching for the value of the regression coefficient(s) that minimizes the total error of the model. There are two main types of linear regression:

• Simple linear regression uses only one independent variable
• Multiple linear regression uses two or more independent variables

Simple Linear Regression: The first dataset contains observations about income (in a range of $15k to $75k) and happiness (rated on a scale of 1 to 10) in an imaginary sample of 500 people. The income values are divided by 10,000 to make the income data match the scale of the happiness scores (so a value of $2 represents $20,000, $3 is $30,000, etc.)
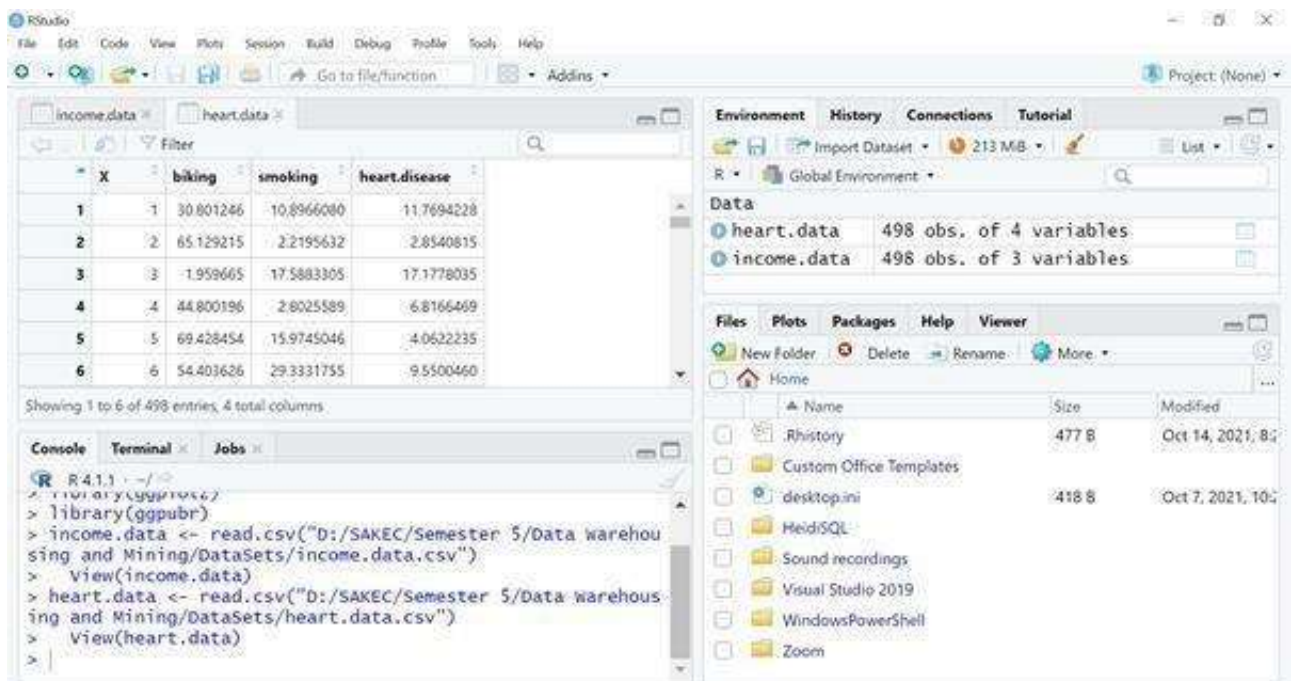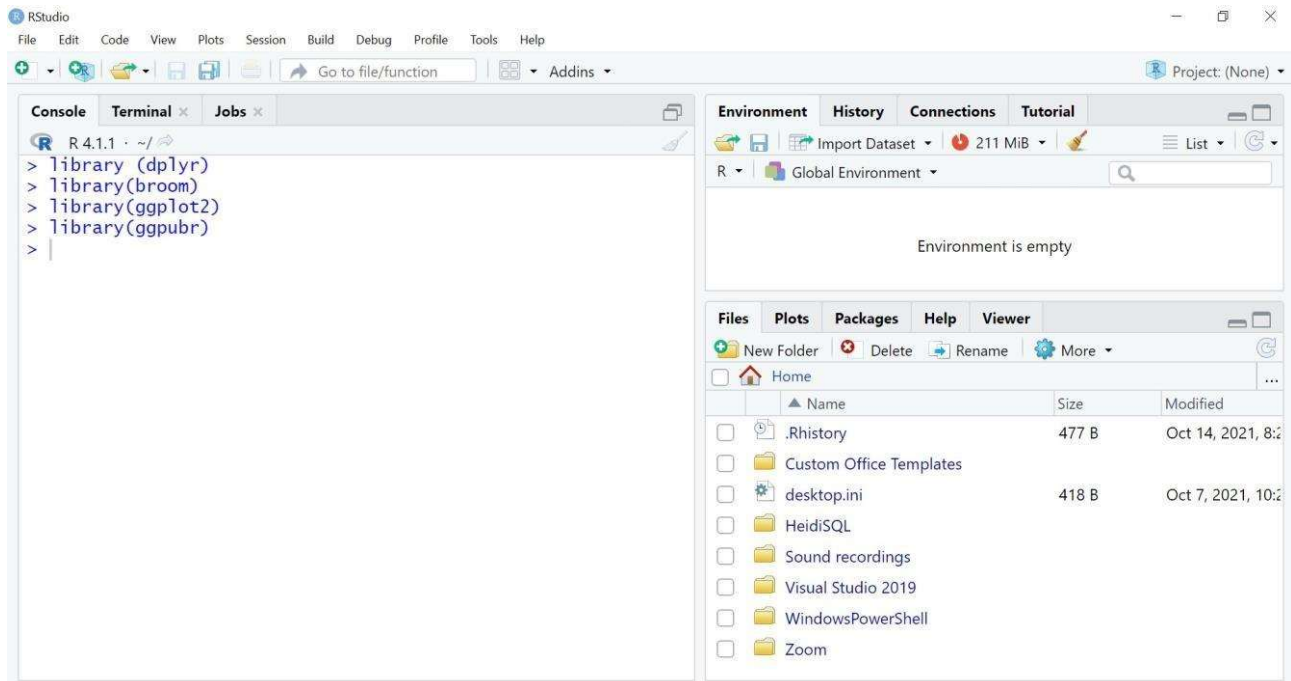
Multiple Linear Regression: The second dataset contains observations on the percentage of people biking to work each day, the percentage of people smoking, and the percentage of people with heart disease in an imaginary sample of 500 towns. Steps to implement Linear Regression in R:

Step 1: Open RStudio and click on File > New File > R Script.
 Then install the packages needed for the analysis, using following code ( only need to do this once):

install.packages("ggplot
2")
install.packages("dplyr"
)
install.packages("broom
")
install.packages("ggpub
r") library(ggplot2)
library(dplyr)
library(broom)
library(ggpubr)

Name:Chavan Dhananjay Rajendra
Div:TE3 Roll No:6
Batch:A

## Step 2: Load the data into R

Follow these four steps for each dataset:

- In RStudio, go to File > Import dataset > From Text (base).
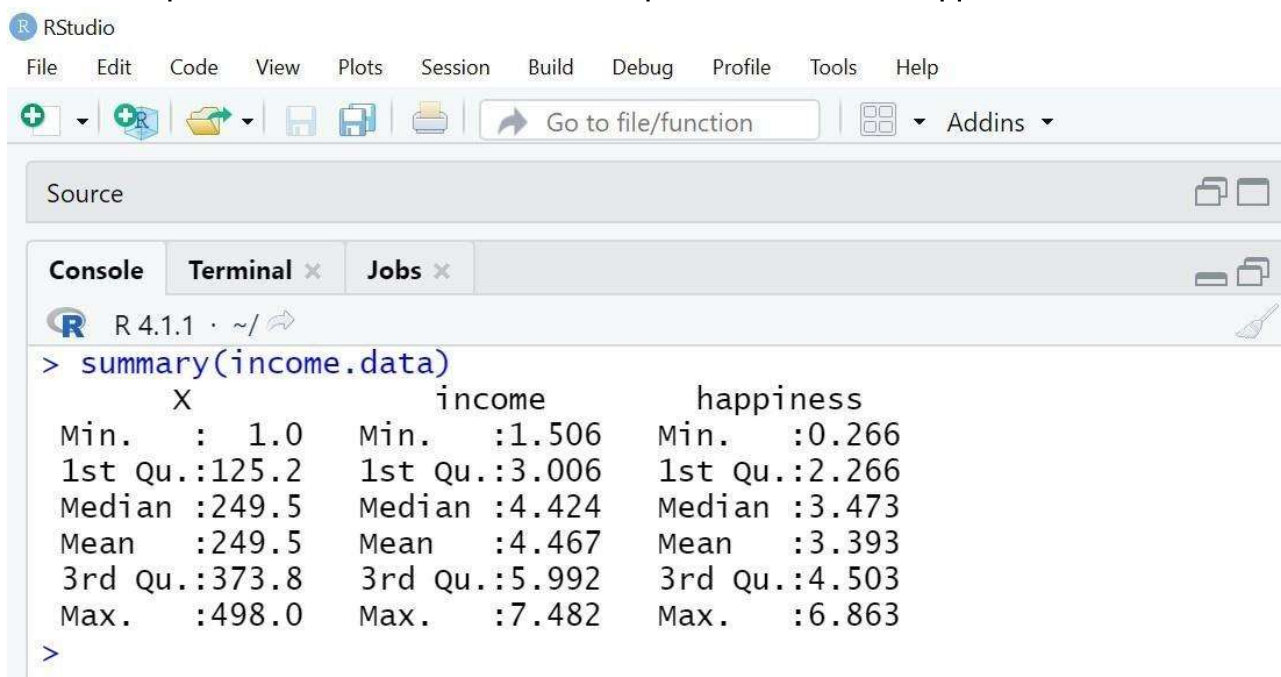
Name:Chavan Dhananjay Rajendra
Div:TE3 Roll No:6
Batch:A

- Choose the data file you have downloaded (income.data or heart. Data), and an Import Dataset window pops up.
- In the Data Frame window, X (index) column and columns listing the data for each of the variables (income and happiness or biking, smoking, and heart.disease) is there.
- Click on the Import button and the file should appear in R Environment tab on the upper right side of the RStudio screen.

Once the data is loaded, check that it has been read in correctly using summary().

1. Simple regression: Because both the variables are quantitative, when this function executed it shows output as a table with a numeric summary of the data. This tells that the minimum, median, mean, and maximum values of the independent variable (income) and dependent variable (happiness):



2. Multiple regression: Again, because the variables are quantitative, running the code produces a numeric summary of the data for the independent variables (smoking and biking) and the dependent variable (heart disease):
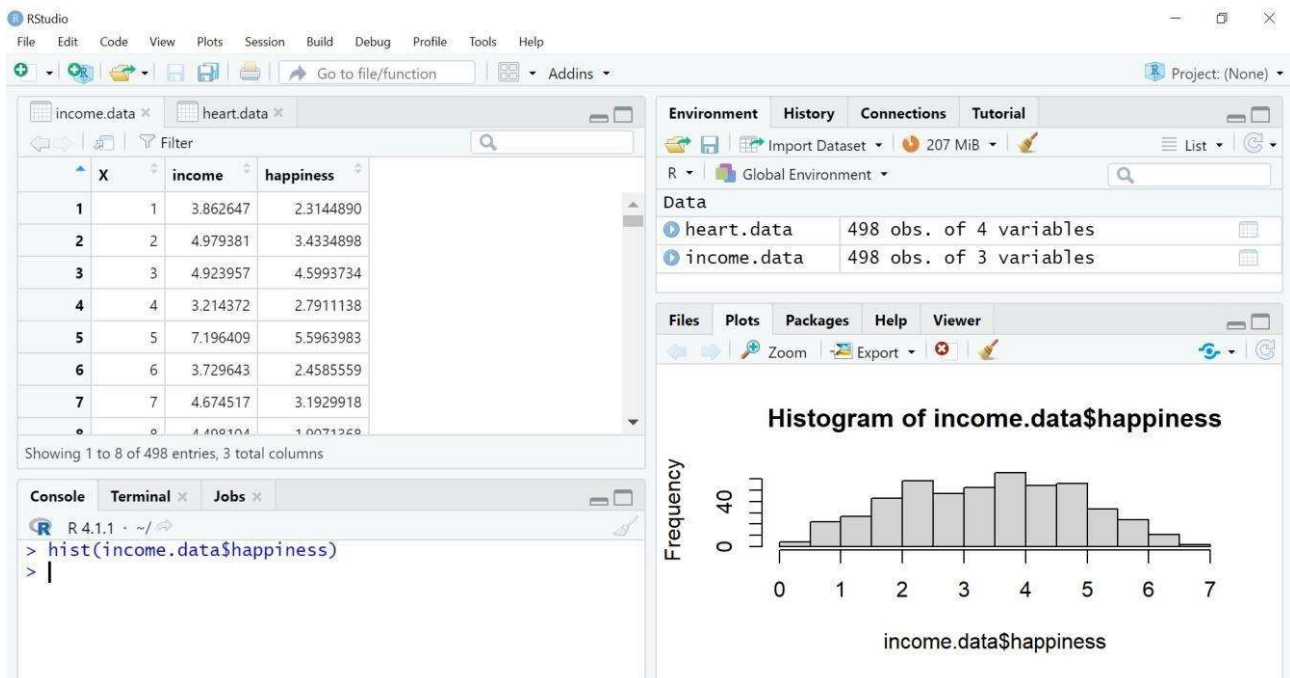
Name:Chavan Dhananjay Rajendra
Div:TE3 Roll No:6
Batch:A



Step 3: To check whether the dependent variable follows a normal distribution, use the hist() function
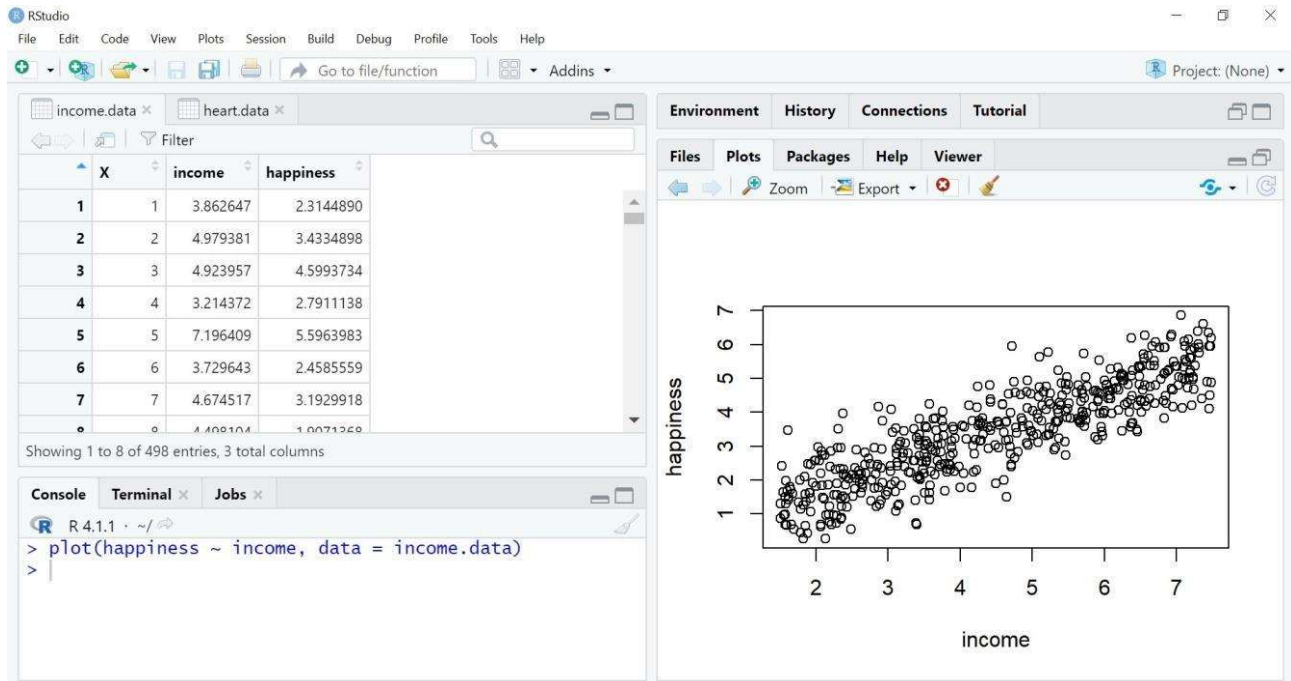
Name:Chavan Dhananjay Rajendra
Div:TE3 Roll No:6
Batch:A

Step 4: The relationship between the independent and dependent variable must be linear. To test this visually with a scatter plot to see if the distribution of data points could be described with a straight line or not.



Step 5: Use the cor() function to test the relationship between independent variables and make sure they aren't too highly correlated.
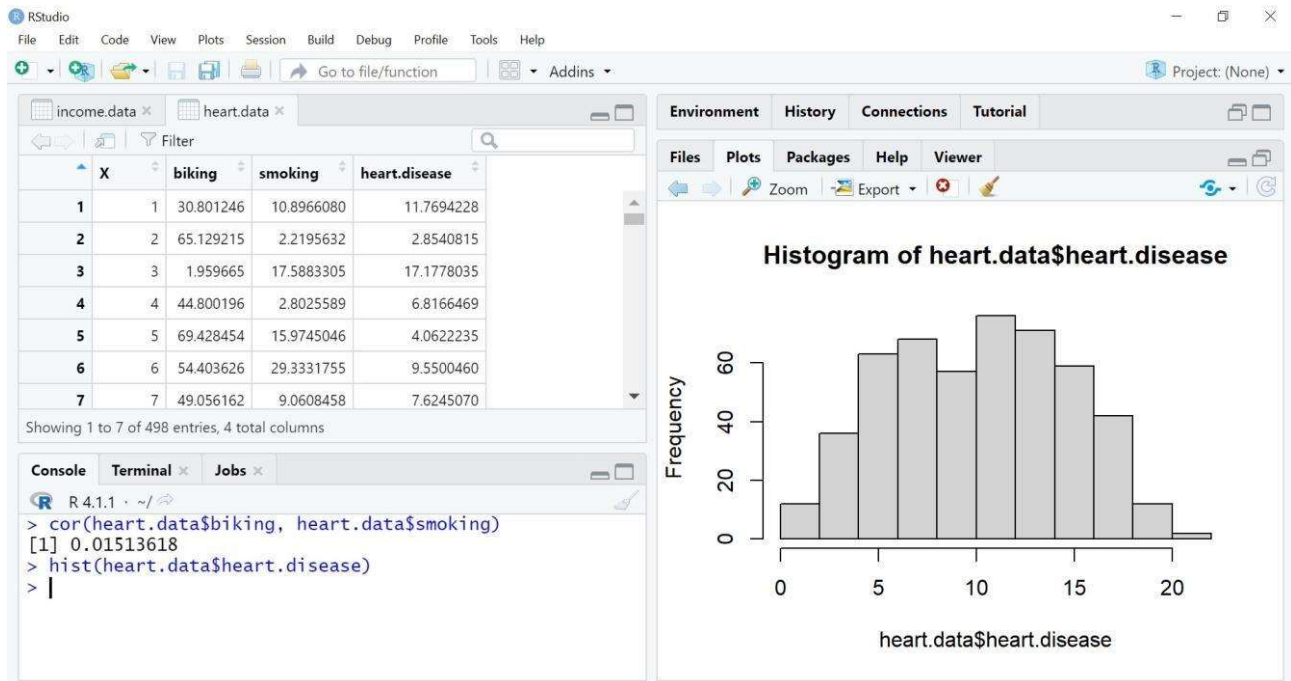
When this code is executed, the output is 0.015. The correlation between biking and smoking is small (0.015 is only a 1.5% correlation), so that include both parameters in our model. Step 6:

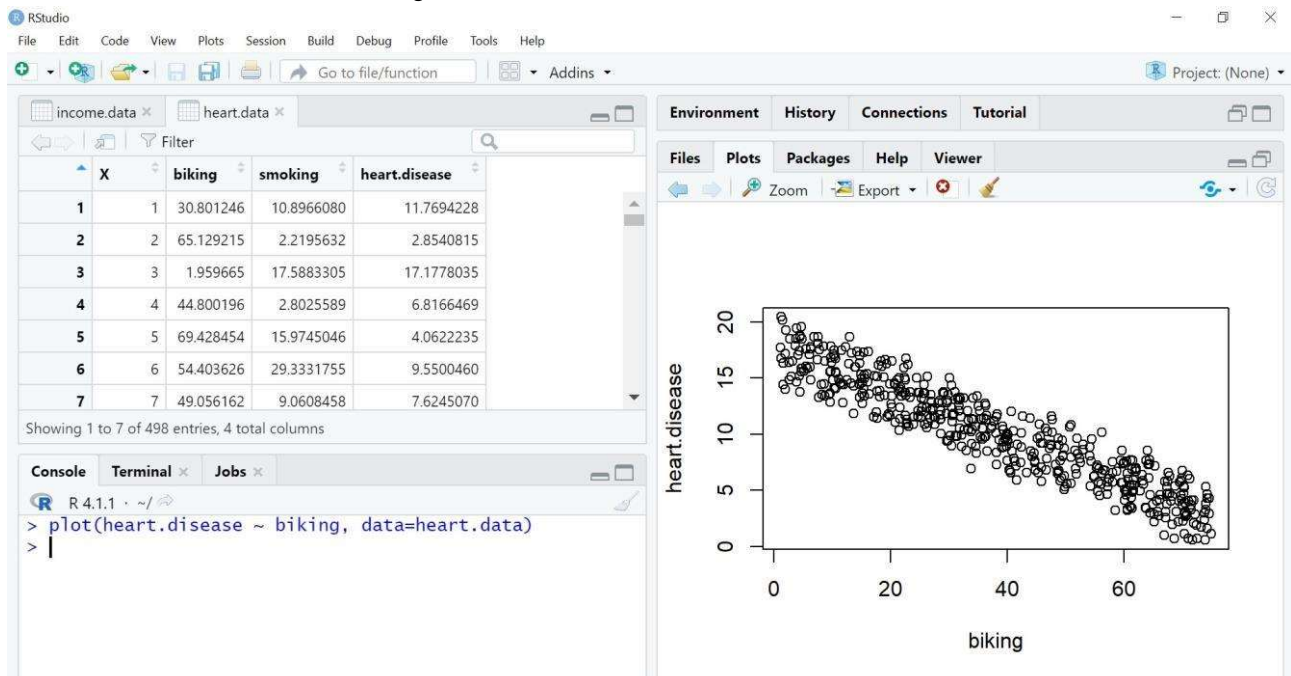Use the hist() function to test whether your dependent variable follows a normal distribution.

**Step 7:** Linearity property is checked using two scatterplots: one for biking and heart disease, and one for smoking and heart disease.
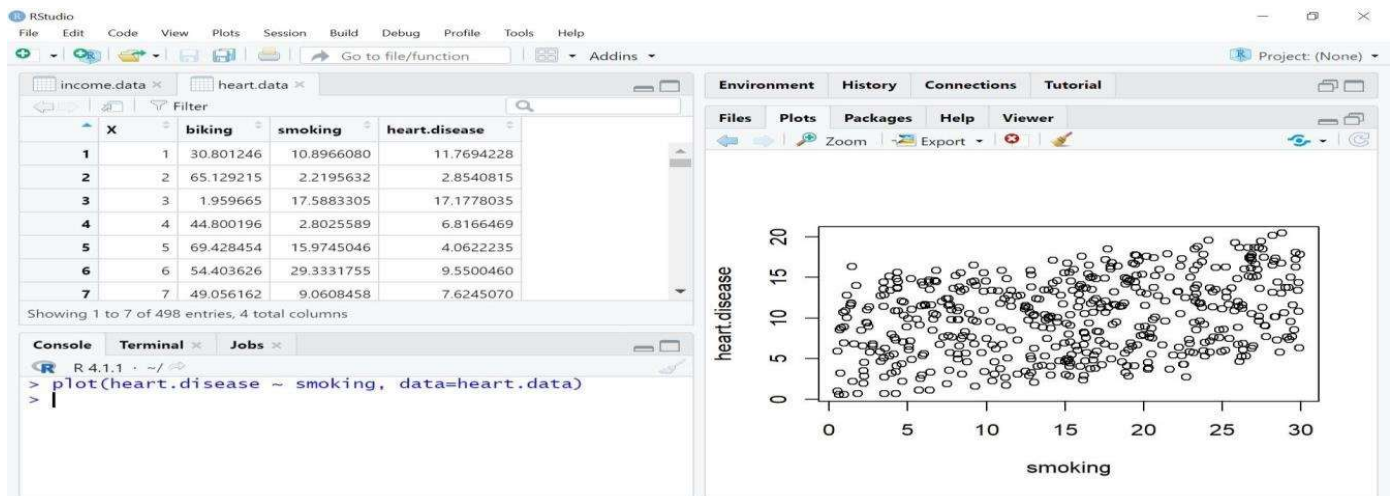
Name:Chavan Dhananjay Rajendra
Div:TE3 Roll No:6
Batch:A



## Step 8: Perform the linear regression analysis

When the data meet the assumptions, perform a linear regression analysis to evaluate the relationship between the independent and dependent variables.

## A. Simple regression: income and happiness



Conclusion:

The above result shows that there is a significant positive relationship between income and happiness ($p$-value < 0.001), with a 0.713-unit (+/- 0.01) increase in happiness for every unit increase in income.