# Data and Artificial Intelligence
# Cyber Shujaa Program

## Week 7 Assignment
### Linear Regression Model

**Student Name:** Faith Jeptoo

**Student ID:** CS-DA02-25005

## Table of Contents

## Introduction

This project focuses on implementing a Simple Linear Regression model to understand and predict a dependent variable based on one independent variable.
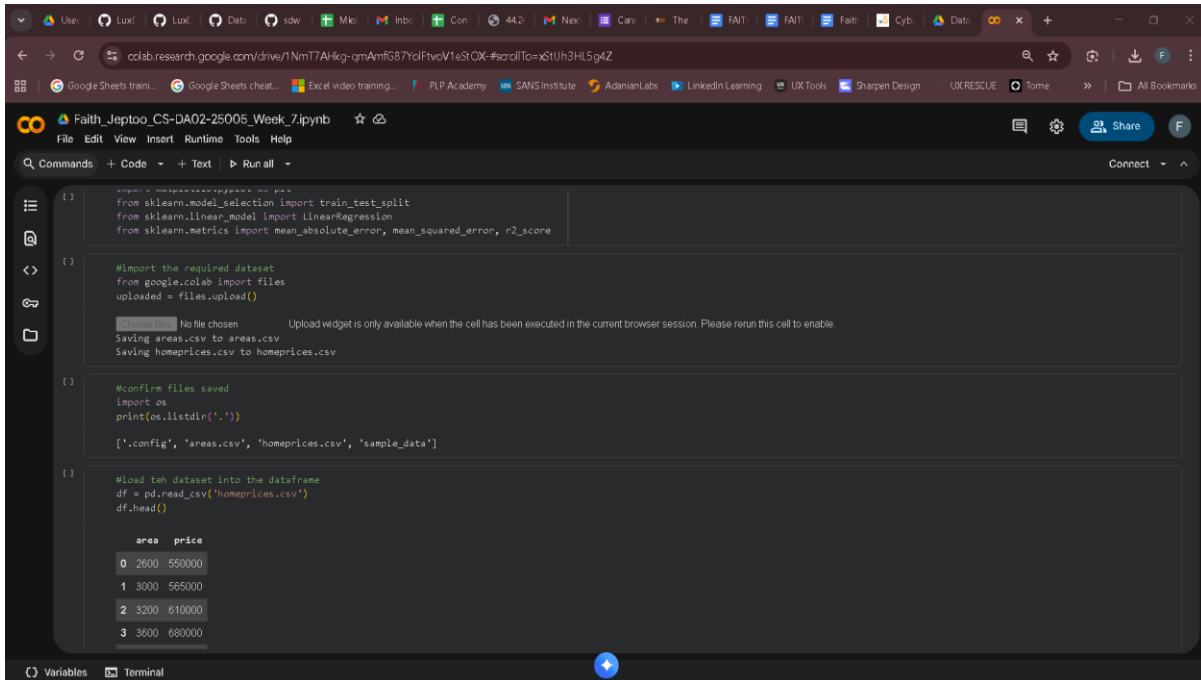
Linear regression is one of the simplest and most widely used statistical algorithms in data science. It helps in establishing a relationship between variables by fitting a straight line to the data.

In this assignment, the primary objective was to explore the dataset, prepare data, train a regression model, evaluate model performance, and visualize results

# Task Completion

## 1 Dataset Exploration

In this section, the dataset was loaded and examined to understand its structure, size, and features. Summary statistics were generated to analyze the distribution of values, and missing values were checked. Visualization through a scatter plot was used to understand the relationship between the independent and dependent variables before applying the regression model



## 2 Data Preparation

The dataset was checked for missing values, and the data was cleaned accordingly. The independent and dependent variables were selected based on the dataset. The dataset was then split into training and testing sets to ensure the model could be evaluated objectively.

## 3. Model Training

The Simple Linear Regression model was trained using the training data. The LinearRegression() class from Scikit-learn was used to fit the model. Once the model was trained, the slope (coefficient) and intercept were extracted to understand the linear equation used for prediction.



## 4. Model Evaluation

To assess the performance of the trained model, several evaluation metrics were used, including Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and the $R^2$ Score. These metrics help determine how accurately the model predicts real values.

## 5. Visualizations

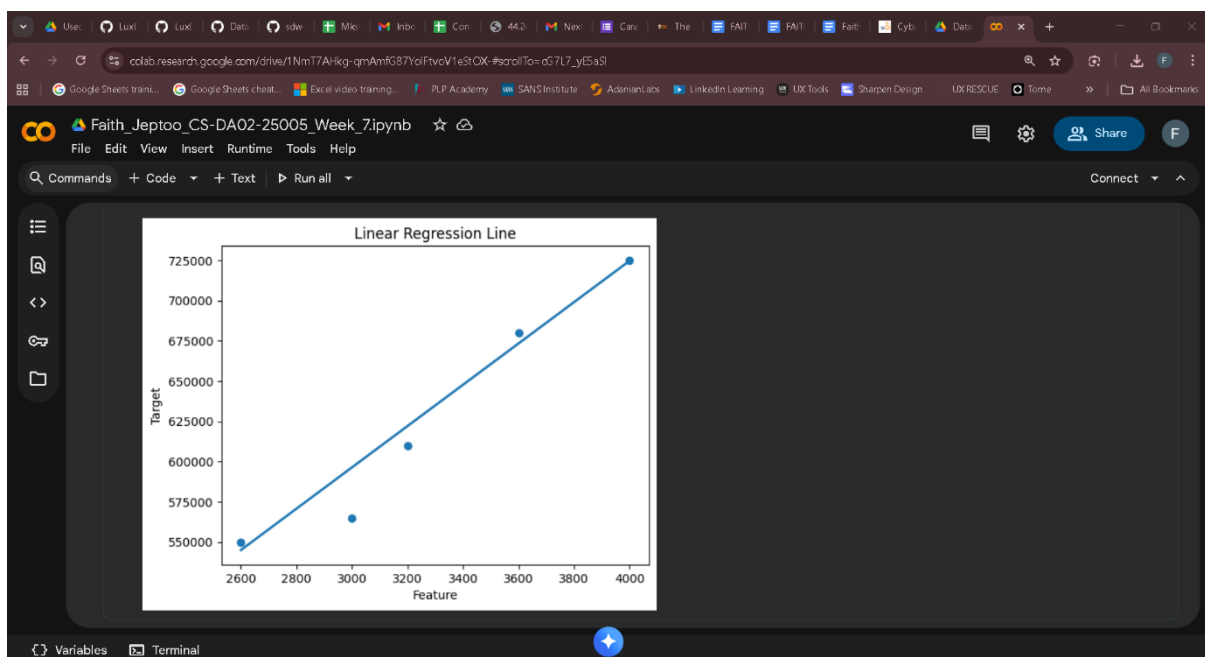Visualizations were created to understand the relationship between variables and validate the performance of the regression model. A scatter plot was used to display the dataset. A regression line was then plotted to show how well the model fits the data visually.

   i.    Scatter plot



   ii.    Linear Regression Line

**6. Conclusion**

This assignment provided hands-on experience in working with a real dataset and applying a simple linear regression model. Through the steps of data exploration, preparation, model training, and evaluation, I gained a deeper understanding of how regression works and how model performance is assessed. Visualization further strengthened the interpretation of results. This knowledge forms a fundamental base for more advanced machine learning models.

**4. Link to Google Colab**

https://colab.research.google.com/drive/1NmT7AHkg-qmAmfG87YolFtvoV1eStOX-?usp=sharing