 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018

Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression

IDE: Google Colab

Theory:

Multiple linear regression (MLR), also known simply as multiple regression, is a statistical technique that uses several explanatory variables to predict the outcome of a response variable. The goal of multiple linear regression is to model the linear relationship between the explanatory (independent) variables and response (dependent) variables. In essence, multiple regression is the extension of ordinary least-squares (OLS) regression because it involves more than one explanatory variable.

Formula and Calculation of Multiple Linear Regression

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon$$

where, for $i = n$ observations:

y_i = dependent variable

x_i = explanatory variables

β_0 = y-intercept (constant term)

β_p = slope coefficients for each explanatory variable


ϵ = the model's error term (also known as the residuals)

Example of How to Use Multiple Linear Regression

As an example, an analyst may want to know how the movement of the market affects the price of ExxonMobil (XOM). In this case, their linear equation will have the value of the S&P 500 index as the independent variable, or predictor, and the price of XOM as the dependent variable.

In reality, multiple factors predict the outcome of an event. The price movement of ExxonMobil, for example, depends on more than just the performance of the overall market. Other predictors such as the price of oil, interest rates, and the price movement of oil futures can affect the price of XOM and stock prices of other oil companies. To understand a relationship in which more than two variables are present, multiple linear regression is used.

Multiple linear regression (MLR) is used to determine a mathematical relationship among several random variables. In other terms, MLR examines how multiple independent variables are related to one dependent variable. Once each of the independent factors has been determined to predict the dependent variable, the information on the multiple variables can be used to create an accurate prediction on the level of effect they have

 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018

on the outcome variable. The model creates a relationship in the form of a straight line (linear) that best approximates all the individual data points.

Methodology:

1. Load the basic libraries and packages
2. Load the dataset
3. Analyse the dataset
4. Normalize the data
5. Pre-process the data
6. Visualize the Data
7. Separate the feature and prediction value columns
8. Write the Hypothesis Function
9. Write the Cost Function
10. Write the Gradient Descent optimization algorithm
11. Apply the training over the dataset to minimize the loss
12. Find the best fit line to the given dataset
13. Observe the cost function vs iterations learning curve

Program (Code):

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

from google.colab import drive #toimport google drive data

drive.mount('/content/drive')


cd /content/drive/"My Drive"

# fetch data
col_name = ['area', 'bedrooms', 'price']
data = pd.read_csv('/content/ex1data2.csv', names=col_name)
data.head(10)

type(dataset)

data.describe()

area_value = dataset.iloc[0:data.shape[0],0:1]
```

 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018

area_value

```
bedroom_value = data.iloc[0:data.shape[0],1:2]
bedroom_value
```

```
price_value = data.iloc[0:47, 2:3]
price_value.head(5)
```

```
# Create 2 subplot, 1 for each variable
fig, axes = plt.subplots(figsize=(12,4),nrows=1,ncols=2)
```

```
axes[0].scatter(area_value,price_value,color="b")
axes[0].set_xlabel("Size (Square Feet)")
axes[0].set_ylabel("Prices")
axes[0].set_title("House prices against size of house")
axes[1].scatter(bedroom_value,price_value,color="r")
axes[1].set_xlabel("Number of bedroom")
axes[1].set_ylabel("Prices")
#axes[1].set_xticks(np.arange(1,6,step=1))
axes[1].set_title("House prices against number of bedroom")
```


```
# Enhance layout
plt.tight_layout()
```

```
# feature normalization
def featureNormalization(x):
    mean = np.mean(x, axis=0)
    std=np.std(x, axis=0)
    X_norm = (x - mean) / std

    return X_norm , mean , std
```

```
data_norm=data.values # converts from dataframe to array
m2=len(data_norm[:,-1]) # last column size
X2=data_norm[:,0:2].reshape(m2,2) # m2 rows and 2 columns
# X2 has all the features i.e. size and bedrooms
X2, mean_X2, std_X2 = featureNormalization(X2) # normalize X2
X2 = np.append(np.ones((m2,1)),X2,axis=1) # append the 1's array in X2 at first
y2=data_norm[:,-1].reshape(m2,1) # price estimate
theta2=np.zeros((3,1))
X2
```

```
# hypothesis
```

 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018

```
def hypothesis(theta_array, x1, x2):
    return theta_array[0] + (theta_array[1]*x1) + (theta_array[2]*x2);

# cost function
def cost_function(theta_array, x1_value, x2_value, y_value, m):
    total_error = 0
    for i in range(m):
        total_error += ((theta_array[0] + theta_array[1]*x1_value[i] + theta_array[2]*x2_value[i]) - y_value[i]) ** 2
    return total_error / 2 * m

# gradient descent
def gradient_descent(theta_array, x1, x2, Y, alpha, m):
    sum_0 = 0
    sum_1 = 0
    sum_2 = 0

    for i in range(m):
        sum_0 += (theta_array[0] + theta_array[1]*x1[i] + theta_array[2]*x2[i]) - Y[i]
        sum_1 += x1[i]*((theta_array[0] + theta_array[1]*x1[i] + theta_array[2]*x2[i]) - Y[i])
        sum_2 += x2[i]*((theta_array[0] + theta_array[1]*x1[i] + theta_array[2]*x2[i]) - Y[i])

    new_theta0 = theta_array[0] - alpha * (sum_0) / m
    new_theta1 = theta_array[1] - alpha * (sum_1) / m
    new_theta2 = theta_array[2] - alpha * (sum_2) / m


    updated_theta_array = [new_theta0, new_theta1, new_theta2]
    return updated_theta_array

# training data
def training(x1_train, x2_train, y_train, alpha, iters):
    m = x1_train.size

    theta_0 = 0
    theta_1 = 0
    theta_2 = 0

    theta_array = [theta_0, theta_1, theta_2]
    cost_function_loss = []

    for i in range(iters):
        theta_array = gradient_descent(theta_array, x1_train, x2_train, y_train, alpha, m)
        cost_function_loss.append(cost_function(theta_array, x1_train, x2_train, y_train, m))
```

 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018

if i % 10 == 0:

```
print('value of theta_0 at iteration %d is: ' % i, theta_array[0])
print('value of theta_1 at iteration %d is: ' % i, theta_array[1])
print('value of theta_2 at iteration %d is: ' % i, theta_array[2], '\n')
ynew = x1_train*theta_array[1] + x2_train*theta_array[2] + theta_array[0]
```

```
x = np.arange(0, len(cost_function_loss), step=1)
plt.plot(x, cost_function_loss, "-b", label="Cost Function Curve")
plt.title("Learning Curve")
plt.xlabel("Number Of Iterations")
plt.ylabel("Cost Function Value")
plt.legend()
plt.show()
print("Cost function values: ",cost_function_loss)
return theta_array
```

price_value.values

alpha = 0.05

iters = 90

```
theta_array = training(X2[:,1:2],X2[:,2:3], y2, alpha, iters)
print("the final value of theta_0 is ",theta_array[0])
print("the final value of theta_1 is ",theta_array[1])
print("the final value of theta_2 is ",theta_array[2])
```

theta_array

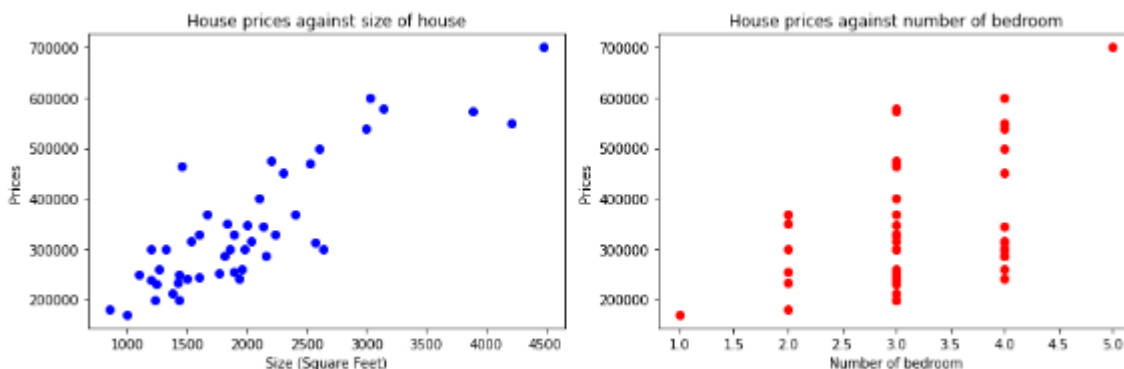
predict = hypothesis(theta_array, 5, 2500)


predict

Results:

To be attached with

- Datapoints scattering (without best fit line)



 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018


b. Data Statistics before Normalization

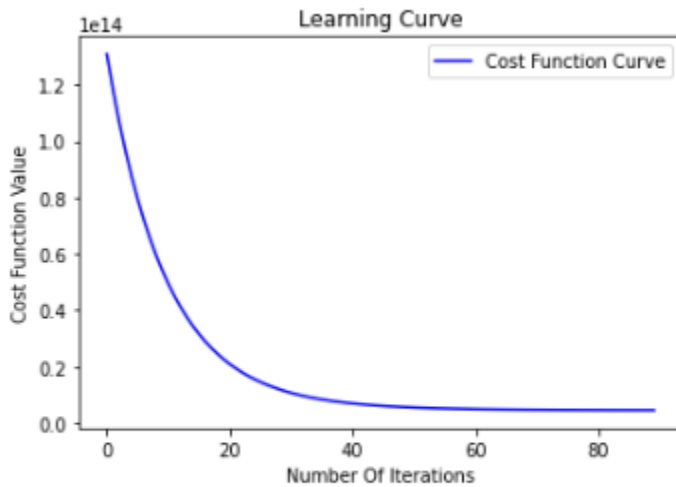
	area	bedrooms	price
count	47.000000	47.000000	47.000000
mean	2000.680851	3.170213	340412.659574
std	794.702354	0.760982	125039.899586
min	852.000000	1.000000	169900.000000
25%	1432.000000	3.000000	249900.000000
50%	1888.000000	3.000000	299900.000000
75%	2269.000000	4.000000	384450.000000
max	4478.000000	5.000000	699900.000000

c. Data Statistics after Normalization

	area	bedrooms	price
count	47.000000	47.000000	47.000000
mean	2000.680851	3.170213	340412.659574
std	794.702354	0.760982	125039.899586
min	852.000000	1.000000	169900.000000
25%	1432.000000	3.000000	249900.000000
50%	1888.000000	3.000000	299900.000000
75%	2269.000000	4.000000	384450.000000
max	4478.000000	5.000000	699900.000000

d. Learning Curve (Cost function vs iterations)


 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018



Cost function values: [array([1.30903637e+14]), array([1.18343329e+14]), array([1.18343329e+14])]
the final value of theta_0 is [337046.53504481]
the final value of theta_1 is [101579.8268961]
the final value of theta_2 is [1220.72747812]

Observation and Result Analysis:

- Nature of the dataset

 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018

b. During Training Process

c. After the training Process


d. Observation over the Learning Curve

Post Lab Exercise:

a. What is the difference between single and multiple variable linear regression

b. What does it mean for a multiple linear regression to be linear?

c. What is the use of Normalization?

 Marwadi University	Marwadi University Faculty of Technology Department of Information and Communication Technology	
Subject: Machine Learning (01CT0607)	Aim: To obtain the best fit line over multiple feature scattered datapoints using Linear Regression	
Experiment No: 03	Date:	Enrollment No:92000133018

- d. Is there any change in the behavior of data before and after normalization? Prove using a toy example.

Post Lab Activity:

Consider any dataset from <https://archive.ics.uci.edu/ml/datasets.php> and perform the multiple variable linear regression analysis over the dataset and obtain the best fit line. Make sure that the dataset is not matching with your classmates. You can also select the dataset from other ML repositories with prior permission from your concerned subject faculty.