 <b>Marwadi University</b>	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

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

**IDE:** Google Colab

### **Theory:**

Linear regression is a method for determining the best linear relationship between two variables  $X$  and  $Y$ . If variables  $X$  and  $Y$  are uncorrelated, it is pointless embarking upon linear regression. However, if a reasonable degree of correlation exists between  $X$  and  $Y$  then linear regression may be a useful means to describe the relationship between the two variables. The usual approach is to use the *least-squares* method, which minimizes the squared difference between the actual data points and a straight line. Let  $[x_i, y_i]$ ,  $i = 1, 2, 3, \dots, N$  be the  $N$  pairs of data values of the variables  $X$  and  $Y$ . The straight-line relating  $X$  and  $Y$  is  $y = mx + c$ , where  $m$  and  $c$  are the gradient and constant values (to be determined) defining the straight line. Thus,  $y(x_i) - y_i$  is the difference between the line and data point  $i$  (see Fig. 1). Taking all the data points, we seek values of  $m$  and  $c$  that minimize the squared difference  $SD$ .

$$\sum_1^N [y(x_i) - y_i]^2$$

This is achieved by calculating the partial derivatives of  $SD$  with respect to  $m$  and  $c$  and finding the pair  $[m, c]$  such that  $SD$  is at a minimum.

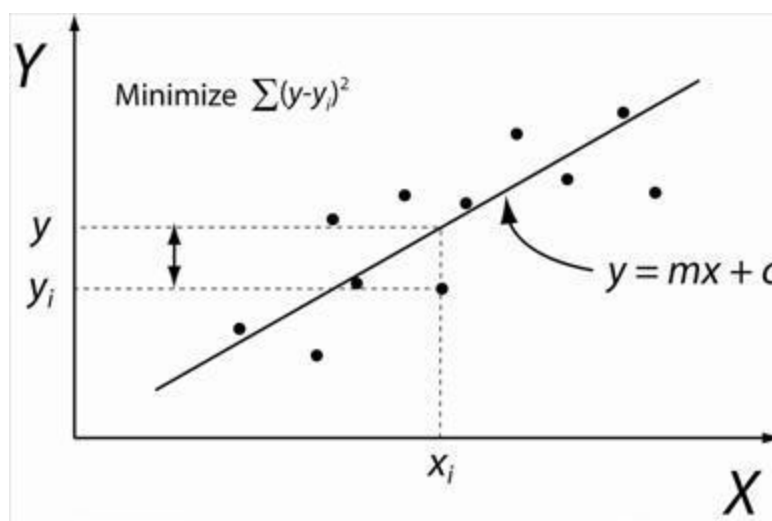



Figure 1: Illustration of Linear Regression. Linear least squares regression, the idea is to find the line  $y = mx + c$  that minimizes the mean squared difference between the line and the data points

 <b>Marwadi University</b>	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

## Batch Gradient Descent:

Gradient Descent is an optimization algorithm used for minimizing the cost function in various machine learning algorithms. It is basically used for updating the parameters of the learning model. Batch gradient descent which processes all the training examples for each iteration of gradient descent. But if the number of training examples is large, then batch gradient descent is computationally very expensive.

Let  $m$  be the number of training examples. Let  $n$  be the number of features.

### Algorithm for batch gradient descent :

Let  $h_{\theta}(x)$  be the hypothesis for linear regression. Then, the cost function is given by:

Let  $\Sigma$  represents the sum of all training examples from  $i=1$  to  $m$ .

$$J_{\text{train}}(\theta) = (1/2m) \Sigma (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

$$\theta_j = \theta_j - (\text{learning rate}/m) * \Sigma (h_{\theta}(x^{(i)}) - y^{(i)})x_j^{(i)}$$

For every  $j = 0 \dots n$

}

Where  $x_j^{(i)}$  Represents the  $j^{\text{th}}$  feature of the  $i^{\text{th}}$  training example. So if  $m$  is very large (e.g. 5 million training samples), then it takes hours or even days to converge to the global minimum. That's why for large datasets, it is not recommended to use batch gradient descent as it slows down the learning.

### Pre Lab Exercise:

a. Explain the meaning of linear regression

---



---



---



---

b. Write three applications of linear regression

---




---



---



---

 <b>Marwadi</b> University	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

c. Write three advantages of linear regression

---



---



---

d. Write three limitations of linear regression

---



---




---

### **Methodology:**

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

### **Program (Code):**

1. Load the basic libraries and packages

 <b>Marwadi</b> University	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>



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

## 2. Load the dataset




```
dataset=pd.read_csv("/content/train.csv")
dataset
```



	x	y
0	24.0	21.549452
1	50.0	47.464463
2	15.0	17.218656
3	38.0	36.586398
4	87.0	87.288984
...	...	...
695	58.0	58.595006
696	93.0	94.625094
697	82.0	88.603770
698	66.0	63.648685
699	97.0	94.975266

700 rows × 2 columns

## 3. Analyse the dataset

 <b>Marwadi</b> University	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

```
[ ] dataset.shape
```


```
⇒ (700, 2)
```

```
[ ] dataset.describe()
```

```
⇒
```

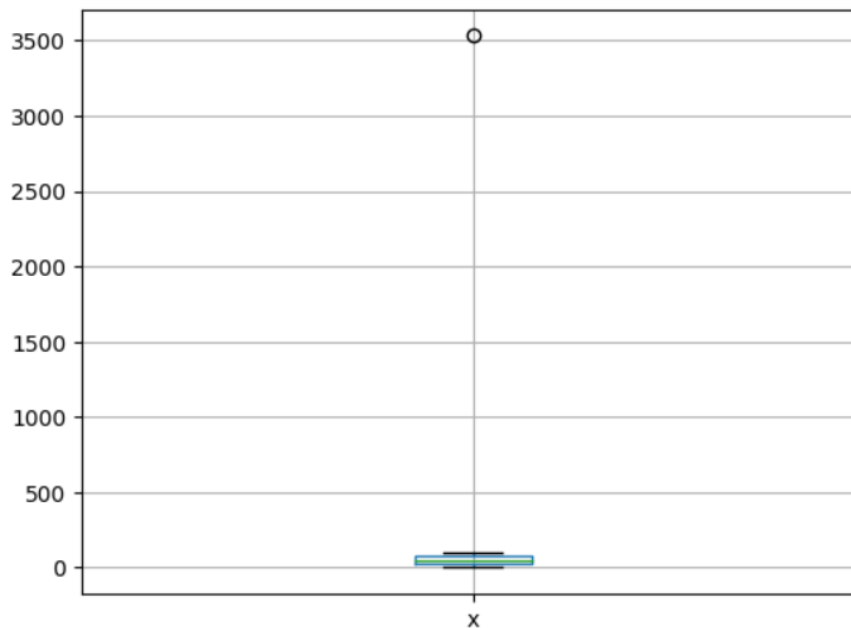
	x	y
count	700.000000	699.000000
mean	54.985939	49.939869
std	134.681703	29.109217
min	0.000000	-3.839981
25%	25.000000	24.929968
50%	49.000000	48.973020
75%	75.000000	74.929911
max	3530.157369	108.871618

#### 4. Pre-process the data

 <b>Marwadi University</b>	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

```
[ ] x_value=dataset.iloc[0:700,0:1]
    y_value=dataset.iloc[0:700,1:2]
    x_value.boxplot(column=['x'])
```

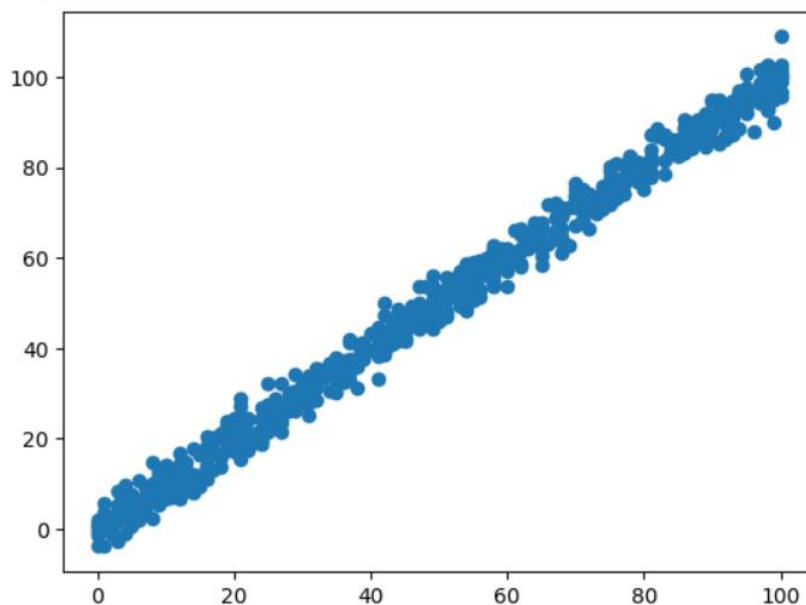
<Axes: >




## 5. Visualize the Data

```
plt.scatter(x_value,y_value)
```

<matplotlib.collections.PathCollection at 0x796eb8128c10>



 <b>Marwadi University</b>	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

## 6. Write the Hypothesis Function

```

7. def hypothesis(theta_array,x):
8.     h=theta_array[0]+theta_array[1]*x
9.     return h

```

## 10. Write the Cost Function

```

11. def costfunction(theta_array,x,y,m):
12.     total_cost=0
13.     for i in range(m):
14.         total_cost+=((theta_array[0]+theta_array[1]*x[i])-y[i])**2
15.     return total_cost/(2*m)

```

## 16. Write the Gradient Descent optimization algorithm

```

def gradient_descent(theta_array,x,y,m,alpha):
    summation_0=0
    summation_1=0
    for i in range(m):
        summation_0+=((theta_array[0]+theta_array[1]*x[i])-y[i])
        summation_1+=((theta_array[0]+theta_array[1]*x[i])-y[i])*x[i]
    new_theta0=theta_array[0]-(summation_0*alpha)/m
    new_theta1=theta_array[1]-(summation_1*alpha)/m
    improvised_theta=[new_theta0,new_theta1]
    print(improvised_theta)
    return improvised_theta

```

## 17. . Apply the training over the dataset to minimize the loss


```

def training(x,y,alpha,epochs):
    theta_0=0
    theta_1=0
    m=x.size
    cost_values=[]
    theta_array=[theta_0,theta_1]
    for i in range(epochs):
        theta_array=gradient_descent(theta_array,x,y,m,alpha)
        loss=costfunction(theta_array,x,y,m)
        cost_values.append(loss)
        y_new=theta_array[0]+theta_array[1]*x

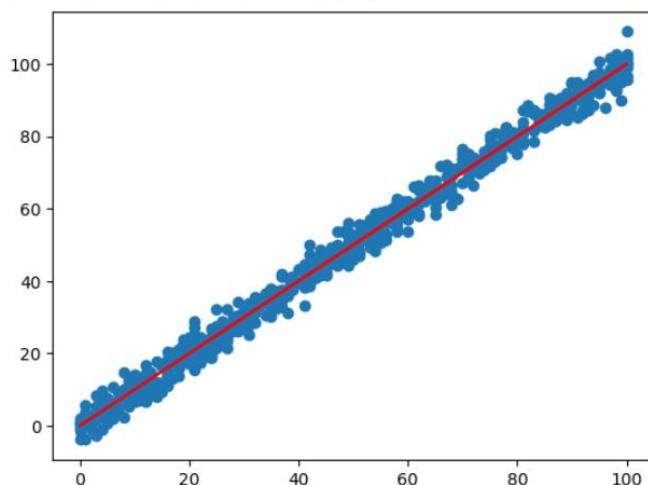
    plt.scatter(x,y)
    plt.plot(x,y_new, 'r')
    plt.show()

    x=np.arange(0,epochs)
    plt.plot(x,cost_values)
    plt.show()

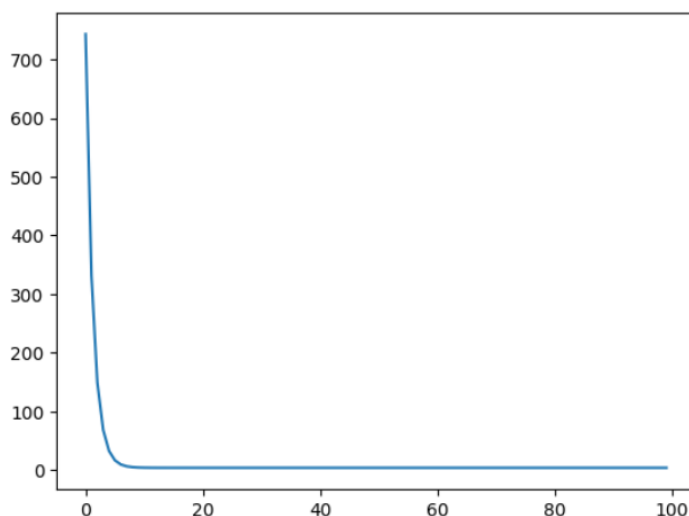
```

 <b>Marwadi University</b>	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

18. Find the best fit line to the given dataset



19. Observe the cost function vs iterations learning curve



### Observation and Result Analysis:

a. Nature of the dataset

---



---




---



---



 <b>Marwadi</b> University	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

b. During Training Process

---

---

---

---

c. After the training Process

---

---

---

---

d. Observation over the Learning Curve

---

---

---

### **Post Lab Exercise:**

a. What are the major assumptions considered in linear regression

---

---

---

b. Why MSE is used instead of MAE for calculating the loss function

---

---


---

c. How can the behaviour of outliers be understood while dealing with the unseen dataset

---

---

---

 <b>Marwadi</b> University	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence(01CT0616)</b>	<b>Aim: To obtain the best fit line over single feature scattered datapoints using Linear Regression</b>	
<b>Experiment No: 02</b>	<b>Date:</b>	<b>Enrollment No: 92200133003</b>

d. Derive the Normal Equation for the Linear Regression.