

Assignment on Multi Regression:

Build an application where it can predict price of a house using multiple variable Linear Regression (use USA Housing dataset from Kaggle). Display all co-efficients and MSE.

```
In [23]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
```

Load the dataset from a CSV file

```
In [24]: df=pd.read_csv(r'C:/Users/3yearb1/Desktop/ML_datasets/USA_Housing.csv')
```

Display the first few rows of the dataset

```
In [25]: display(df.head())
```

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
0	79545.458574	5.682861	7.009188	4.09	23086.800503	1.059034e+06	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	79248.642455	6.002900	6.730821	3.09	40173.072174	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.067179	5.865890	8.512727	5.13	36882.159400	1.058988e+06	9127 Elizabeth Stravenue\nDanieltown, WI 06482...
3	63345.240046	7.188236	5.586729	3.26	34310.242831	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.197226	5.040555	7.839388	4.23	26354.109472	6.309435e+05	USNS Raymond\nFPO AE 09386

Display information about the dataset

In [26]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 7 columns):
Avg. Area Income      5000 non-null float64
Avg. Area House Age   5000 non-null float64
Avg. Area Number of Rooms  5000 non-null float64
Avg. Area Number of Bedrooms  5000 non-null float64
Area Population        5000 non-null float64
Price                  5000 non-null float64
Address                5000 non-null object
dtypes: float64(6), object(1)
memory usage: 273.5+ KB
```

Display summary statistics

In [27]: `df.describe()`

Out[27]:

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price
count	5000.000000	5000.000000	5000.000000	5000.000000	5000.000000	5.000000e+03
mean	68583.108984	5.977222	6.987792	3.981330	36163.516039	1.232073e+06
std	10657.991214	0.991456	1.005833	1.234137	9925.650114	3.531176e+05
min	17796.631190	2.644304	3.236194	2.000000	172.610686	1.593866e+04
25%	61480.562388	5.322283	6.299250	3.140000	29403.928702	9.975771e+05
50%	68804.286404	5.970429	7.002902	4.050000	36199.406689	1.232669e+06
75%	75783.338666	6.650808	7.665871	4.490000	42861.290769	1.471210e+06
max	107701.748378	9.519088	10.759588	6.500000	69621.713378	2.469066e+06

Split the dataset into input variables (X) and output variable (y)

In [28]: `x=df[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms'], 'Area Population']`
`y=df['Price']`

Split the data into training and testing sets (80% training, 20% testing)

In [29]: `x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_s`

Initialize and train a linear regression model

In [30]: `model=LinearRegression()`
`model.fit(x_train,y_train)`

Out[30]: `LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)`

Make predictions on the test set

```
In [31]: y_pred=model.predict(x_test)
```

Calculate Mean Squared Error (MSE) and R squared score

```
In [32]: mse=mean_squared_error(y_test,y_pred)
r2=r2_score(y_test,y_pred)
print(f"Mean squared error:{mse}")
print(f"R-squared score: {r2}")
```

Mean squared error:10073721633.872011

R-squared score: 0.9181214278738137

Visualize the relationship between actual and predicted prices

```
In [ ]:
```