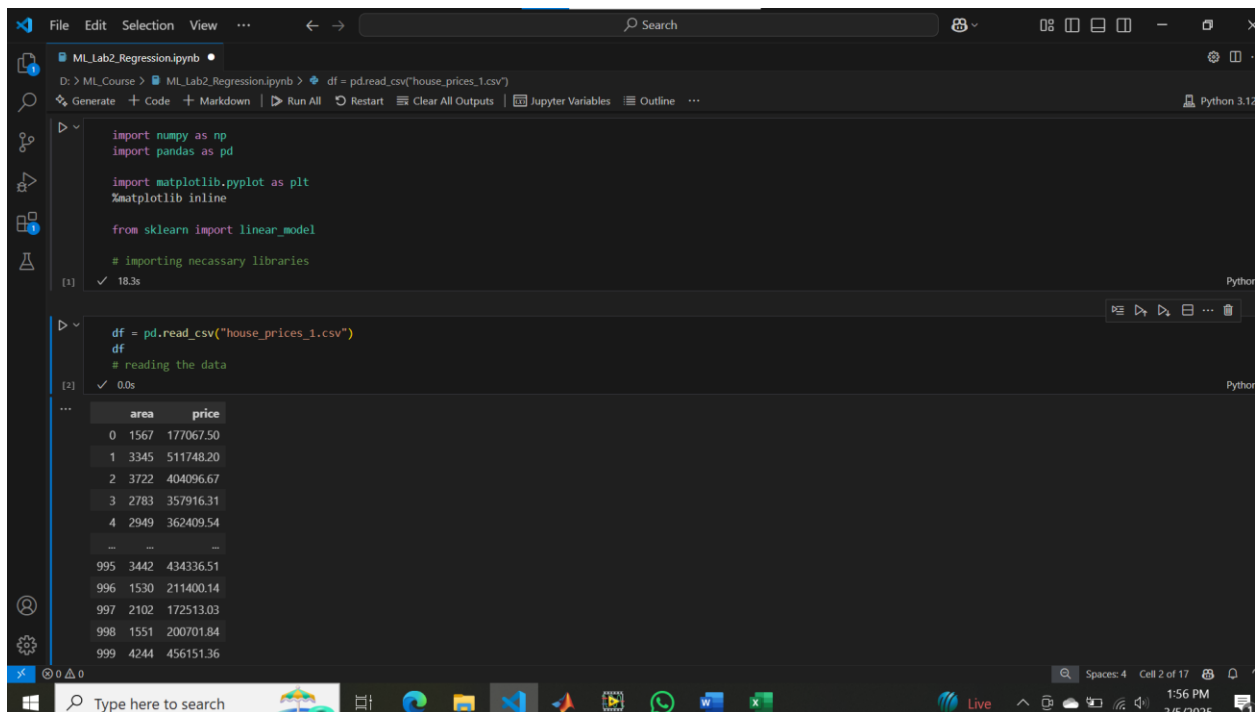# ML_Lab2_Regression

**Student Name: Abdul Malik**

**Student Number: 2211011098**

**Objective:**

Make a Linear Regression model that predicts the house price given a specific feature i.e. area in square feet.



The above two cells of jupyter notebook show the necessary libraries which are firstly the famous ones, numpy, panda, matplotlib (for plotting the data) and scikit learn (famous for having packages for linear regression, logistic regression, classification and clustering etc) ([scikit-learn: machine learning in Python — scikit-learn 1.6.1 documentation](#)).

Panda's **read_csv** reads the data from the saved **dataset house_prices_1.csv** and then prints the elements of csv file.

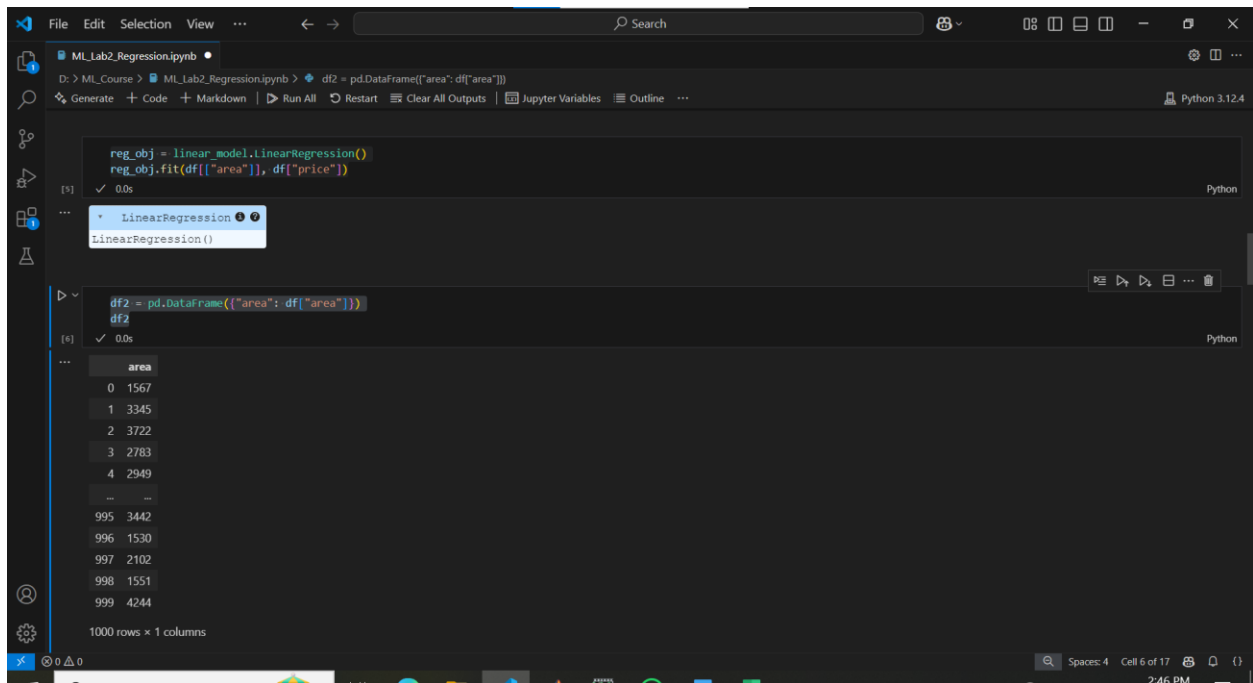Now the datatypes of the elements inside csv (df variable defined) are shown using **dtypes** from Pandas.



Here we scatter plot the elements of datasets which are area (square feet) and price ($) using **plt** function from **matplotlib.**

ML_Lab2_Regression.ipynb  ●

D: > ML_Course > 🔷 ML_Lab2_Regression.ipynb > 🔷 df2 = pd.DataFrame({"area": df["area"]})

✦ Generate  + Code  + Markdown  | ▷ Run All  ↺ Restart  ☰ Clear All Outputs  | 🔲 Jupyter Variables  ☰ Outline  ...          🐍 Python 3.12.4

```python
reg_obj = linear_model.LinearRegression()
reg_obj.fit(df[["area"]], df["price"])
```
[5]  ✓ 0.0s                                                                                                           Python

```
  ▾  LinearRegression ❶ ❷
LinearRegression()
```

```python
df2 = pd.DataFrame({"area": df["area"]})
df2
```
[6]  ✓ 0.0s                                                                                                           Python

|     | area |
| --- | ---- |
| 0   | 1567 |
| 1   | 3345 |
| 2   | 3722 |
| 3   | 2783 |
| 4   | 2949 |
| ... | ...  |
| 995 | 3442 |
| 996 | 1530 |
| 997 | 2102 |
| 998 | 1551 |
| 999 | 4244 |

1000 rows × 1 columns

Here we load the features (area in our case) to our linear regression model form scikit learn.

```python
predicted_arr = reg_obj.predict(df2)

print(predicted_arr[:10], end=", ")
print("...", end=", ")
print(predicted_arr[-10:])
```

```
[212322.90487091 384074.21961441 420491.68061458 329786.12125342
 345821.39585827 405712.18052698 266417.80715232 360600.89594587
 484342.98491461 276657.19936988], ..., [378857.92546585 212709.29703006 416724.35706284 393251.0333943
 468597.50442913 393444.22947387 208748.77739874 264002.85615762
 210777.33623429 470915.85738405]
```

```python
m = reg_obj.coef_  # ouputs the value of slope(m)
m
```

```
array([96.59803979])
```

```python
c = reg_obj.intercept_  # outputs the value of intercept(c)
c
```

```
60953.776522719534
```

```python
x = df2.to_numpy(df["area"])
```

```python
x = df2.to_numpy(df["area"])
```

```python
# y = w * x + b

using_formula_arr = np.empty(1000)
for i in range(len(x)):
    using_formula_arr[i] = w * x[i] + b
#    print(m * x[i] + c, end=", ")

print(using_formula_arr[:10], end=", ")
print("...", end=", ")
print(using_formula_arr[-10:])
```

```
[212322.90487091 384074.21961441 420491.68061458 329786.12125342
 345821.39585827 405712.18052698 266417.80715232 360600.89594587
 484342.98491461 276657.19936988], ..., [378857.92546585 212709.29703006 416724.35706284 393251.0333943
 468597.50442913 393444.22947387 208748.77739874 264002.85615762
 210777.33623429 470915.85738405]
C:\Users\DELL\AppData\Local\Temp\ipykernel_17788\3950902830.py:5: DeprecationWarning: Conversion of an array with ndim > 0 to a scalar is deprecated, and will error in future. Ensure yo
    using_formula_arr[i] = w * x[i] + b
```

```python
np.array_equal(predicted_arr, using_formula_arr)
```

```
True
```

Using the prediction equation y' = w*x + b, we calculate prediction inside a for loop by running it from i to number of features which is **len(x)** and saving each iteration inside using_formula_arr[i] array.

Now we pass the using_formula_arr[i] and predicted_arr (predicted value) to a numpy array and display the predicted_prices ($) together with our feature (area) and price ($).
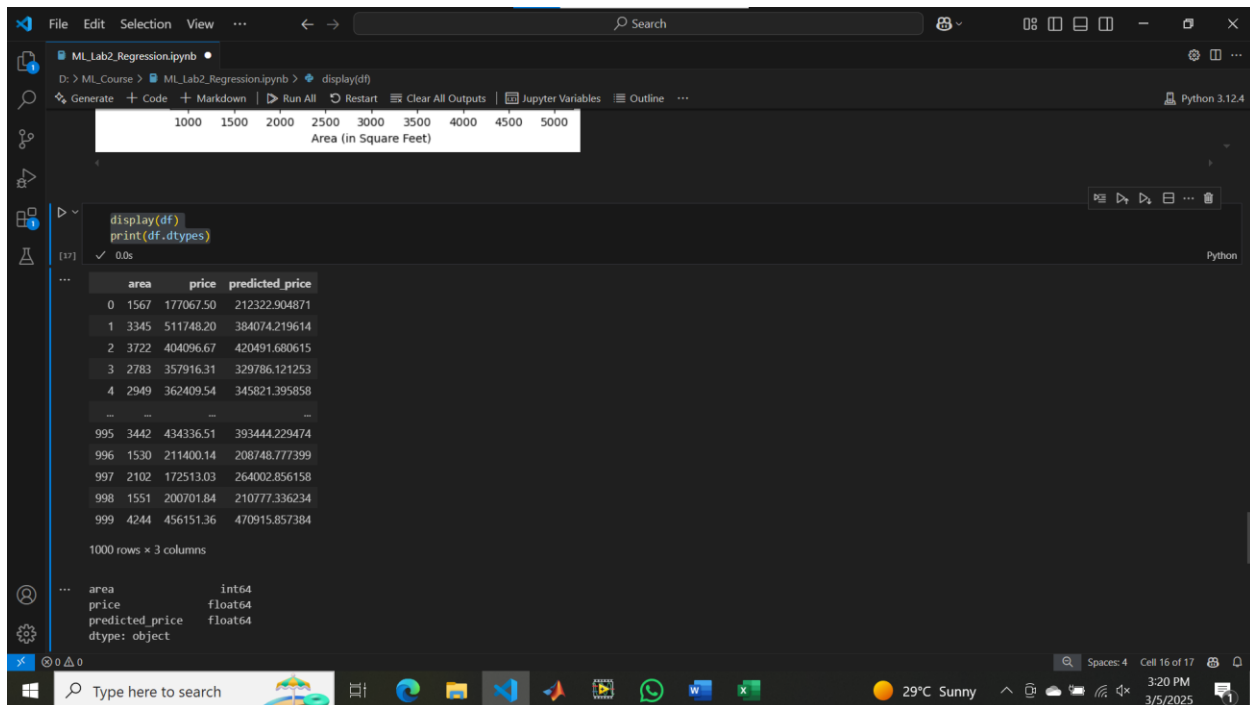


As shown above, plotting the linear regression best fit line. Why we choose the lime-colored line here? It is because we are calculating errors which is the vertical distance between our actual values and predicted values. We squared the errors that we found, taking their squares and
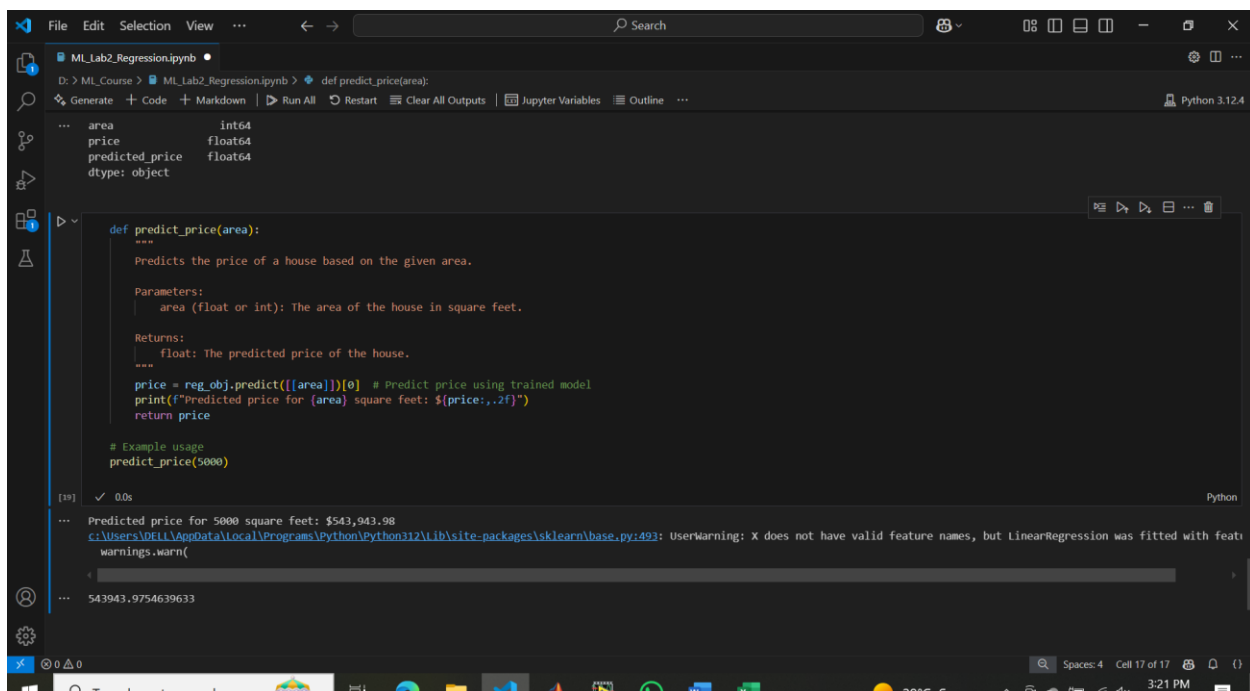
summing them up and making sure we minimize that error. So my lime-colored line actually represents the minimum sum value.

Now in a single scatter plot, we combine the regression line and the dataset values that we have to get the prediction model.



Displaying predicted price.



Here we are creating a function **predict_price** and giving it the **area** as a parameter. Here we can pass a certain value of our choice for the area that we want to predict the price for. As we can

see when an area of **5000 square feet** is passed to the function, it predicts the price of **543943.98$.** It can be verified with the table given above as well.

**Conclusion:**

Our model gave us good results at the end. We implemented the simple linear regression model and got pretty close values to our actual values with almost 82% accuracy. If we train it with more datasets and different features we will get a more accurate result.