

Business Problem

```
In [1]: import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt
```

Load Dataset

```
In [2]: df = pd.read_csv(r'D:\Samson - All Data\Naresh IT Institute\New folder\student_info.csv')
```

```
In [3]: df
```

Out[3]:

	study_hours	student_marks
0	6.83	78.50
1	6.56	76.74
2	NaN	78.68
3	5.67	71.82
4	8.67	84.19
...
195	7.53	81.67
196	8.56	84.68
197	8.94	86.75
198	6.60	78.05
199	8.35	83.50

200 rows × 2 columns

In [4]:

`df.head()`

Out[4]:

	study_hours	student_marks
0	6.83	78.50
1	6.56	76.74
2	NaN	78.68
3	5.67	71.82
4	8.67	84.19

In [6]:

`df.shape`

Out[6]: (200, 2)

In [10]: `df.tail()`

Out[10]:

	study_hours	student_marks
195	7.53	81.67
196	8.56	84.68
197	8.94	86.75
198	6.60	78.05
199	8.35	83.50

Discover and visualize the data to gain insights

In [7]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype  
---  -
0   study_hours     195 non-null   float64
1   student_marks   200 non-null   float64
dtypes: float64(2)
memory usage: 3.3 KB
```

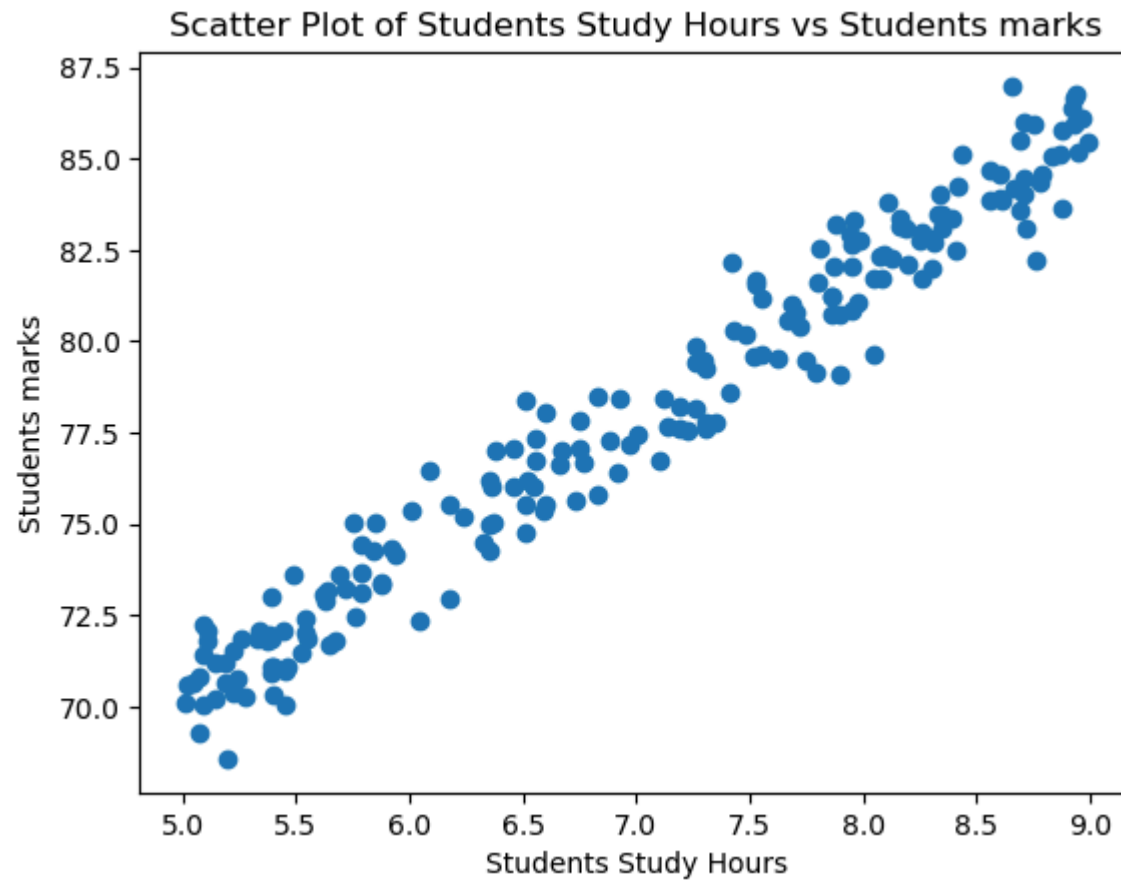
Descriptive statistics

In [11]: `df.describe()`

Out[11]:

	study_hours	student_marks
count	195.000000	200.000000
mean	6.995949	77.93375
std	1.253060	4.92570
min	5.010000	68.57000
25%	5.775000	73.38500
50%	7.120000	77.71000
75%	8.085000	82.32000
max	8.990000	86.99000

```
In [13]: plt.scatter(x = df.study_hours, y = df.student_marks)
plt.xlabel("Students Study Hours")
plt.ylabel("Students marks")
plt.title("Scatter Plot of Students Study Hours vs Students marks")
plt.show()
```



Prepare the data for Machine Learning algorithms

```
In [14]: # Data Cleaning
```

```
In [15]: df
```

Out[15]:

	study_hours	student_marks
0	6.83	78.50
1	6.56	76.74
2	NaN	78.68
3	5.67	71.82
4	8.67	84.19
...
195	7.53	81.67
196	8.56	84.68
197	8.94	86.75
198	6.60	78.05
199	8.35	83.50

200 rows × 2 columns

In [18]: `df.isnull().sum()`Out[18]:

```
study_hours    5
student_marks  0
dtype: int64
```

In [20]: `df.mean()`Out[20]:

```
study_hours    6.995949
student_marks  77.933750
dtype: float64
```

In [21]: `df2 = df.fillna(df.mean())`In [22]: `df2.isnull().sum()`

```
Out[22]: study_hours    0
         student_marks  0
         dtype: int64
```

```
In [23]: df.head()
```

```
Out[23]:
```

	study_hours	student_marks
0	6.83	78.50
1	6.56	76.74
2	NaN	78.68
3	5.67	71.82
4	8.67	84.19

```
In [24]: df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   study_hours     200 non-null   float64
1   student_marks   200 non-null   float64
dtypes: float64(2)
memory usage: 3.3 KB
```

```
In [25]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   study_hours     195 non-null   float64
1   student_marks   200 non-null   float64
dtypes: float64(2)
memory usage: 3.3 KB
```

```
In [26]: # split dataset
```

```
In [27]: x = df2.drop("student_marks", axis = "columns")
y = df2.drop("study_hours", axis = "columns")
print("shape of x = ", x.shape)
print("shape of y = ", y.shape)
```

shape of x = (200, 1)

shape of y = (200, 1)

```
In [28]: x
```

```
Out[28]:
```

	study_hours
--	-------------

0	6.830000
---	----------

1	6.560000
---	----------

2	6.995949
---	----------

3	5.670000
---	----------

4	8.670000
---	----------

...	...
-----	-----

195	7.530000
-----	----------

196	8.560000
-----	----------

197	8.940000
-----	----------

198	6.600000
-----	----------

199	8.350000
-----	----------

200 rows × 1 columns

```
In [29]: y
```


Out[29]:

student_marks	
0	78.50
1	76.74
2	78.68
3	71.82
4	84.19
...	...
195	81.67
196	84.68
197	86.75
198	78.05
199	83.50

200 rows × 1 columns

```
In [31]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state=0)

print("shape of x_train = ", x_train.shape)
print("shape of y_train = ", y_train.shape)
print("shape of x_test = ", x_test.shape)
print("shape of y_test = ", y_test.shape)
```

```
shape of x_train = (160, 1)
shape of y_train = (160, 1)
shape of x_test = (40, 1)
shape of y_test = (40, 1)
```

```
In [32]: x_train
```

Out[32]:

study_hours	
134	6.51
66	7.86
26	6.51
113	7.95
168	7.95
...	...
67	8.26
192	8.71
117	8.83
47	5.01
172	7.35

160 rows × 1 columns

In [35]: `x_test.shape`Out[35]: `(40, 1)`In [33]: `y_train`

Out[33]:

student_marks	
134	78.39
66	81.25
26	74.75
113	80.86
168	82.68
...	...
67	81.70
192	84.03
117	85.04
47	70.11
172	77.78

160 rows × 1 columns

In [36]: `y_train.shape`

Out[36]: (160, 1)

In [37]: `y_test`

Out[37]:

student_marks	
18	82.50
170	71.18
107	73.25
98	83.64
177	73.64
182	86.99
5	81.18
146	82.75
12	79.50
152	81.70
61	79.41
125	85.95
180	77.19
154	78.45
80	84.00
7	85.46
33	84.35
130	73.19
37	78.21
74	77.59
183	83.87
145	85.15

student_marks	
45	72.96
159	80.72
60	73.61
123	79.53
179	78.17
185	79.63
122	76.83
44	82.38
16	76.04
55	85.48
150	71.87
111	75.04
22	70.67
189	79.87
129	74.49
4	84.19
83	75.36
106	72.10

In [38]: x_test

Out[38]:

	study_hours
18	8.410000
170	5.190000
107	5.720000
98	8.880000
177	5.790000
182	8.660000
5	7.550000
146	7.990000
12	7.750000
152	8.080000
61	7.260000
125	8.750000
180	6.970000
154	6.930000
80	8.340000
7	8.990000
33	8.780000
130	5.640000
37	7.190000
74	7.310000
183	8.610000
145	8.950000

study_hours	
45	6.180000
159	7.860000
60	5.490000
123	7.620000
179	7.260000
185	8.050000
122	6.995949
44	8.090000
16	6.360000
55	8.690000
150	5.390000
111	6.370000
22	5.050000
189	7.260000
129	6.330000
4	8.670000
83	6.010000
106	5.340000

Select a model and train it

```
In [40]: #  $y = m * x + c$   
from sklearn.linear_model import LinearRegression
```

```
lr = LinearRegression()
```

```
In [41]: lr
```

```
Out[41]:
```

LinearRegression ⓘ ?
LinearRegression()

In a Jupyter environment, please return this cell to show the HTML representation or trust the notebook. On GitHub, the HTML representation is unable to render please try loading this page with nbviewer.org. LinearRegression LinearRegression()

```
In [44]: lr.fit(x_train,y_train)
```

```
Out[44]:
```

LinearRegression ⓘ ?
LinearRegression()

- In a Jupyter environment, please return this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render please try loading this page with nbviewer.org. LinearRegression LinearRegression()

```
In [45]: lr.coef_ #m
```

```
Out[45]: array([[3.93037294]])
```

```
In [46]: lr.intercept_ #c
```

```
Out[46]: array([50.45063632])
```

```
In [47]: m = 3.93  
c = 50.44  
y = m * 10 + c  
y
```


Out[47]: 89.74000000000001

```
In [48]: m = 3.93  
c = 50.44  
y = m * 11 + c  
y
```

Out[48]: 93.67

```
In [49]: lr.predict([[11]]).round(2)
```

C:\Users\samua\anaconda3\Lib\site-packages\sklearn\utils\validation.py:2739: UserWarning: X does not have valid feature names, but LinearRegression was fitted with feature names
warnings.warn(

Out[49]: array([[93.68]])

```
In [50]: y_pred = lr.predict(x_test)  
y_pred
```

```
Out[50]: array([[83.50507271],
               [70.84927186],
               [72.93236952],
               [85.35234799],
               [73.20749562],
               [84.48766595],
               [80.12495199],
               [81.85431608],
               [80.91102657],
               [82.20804964],
               [78.98514384],
               [84.84139951],
               [77.84533568],
               [77.68812077],
               [83.22994661],
               [85.78468901],
               [84.9593107 ],
               [72.61793968],
               [78.71001773],
               [79.18166248],
               [84.2911473 ],
               [85.6274741 ],
               [74.74034107],
               [81.3433676 ],
               [72.02838374],
               [80.40007809],
               [78.98514384],
               [82.09013845],
               [77.94732382],
               [82.24735337],
               [75.44780819],
               [84.60557713],
               [71.63534645],
               [75.48711192],
               [70.29901965],
               [78.98514384],
               [75.32989701],
               [84.52696967],
               [74.07217767],
               [71.4388278 ]])
```

```
In [51]: pd.DataFrame(np.c_[x_test, y_test, y_pred], columns = ["study_hours", "student_marks_original", "student_marks_predicted"])
```

Out[51]:

	study_hours	student_marks_original	student_marks_predicted
0	8.410000	82.50	83.505073
1	5.190000	71.18	70.849272
2	5.720000	73.25	72.932370
3	8.880000	83.64	85.352348
4	5.790000	73.64	73.207496
5	8.660000	86.99	84.487666
6	7.550000	81.18	80.124952
7	7.990000	82.75	81.854316
8	7.750000	79.50	80.911027
9	8.080000	81.70	82.208050
10	7.260000	79.41	78.985144
11	8.750000	85.95	84.841400
12	6.970000	77.19	77.845336
13	6.930000	78.45	77.688121
14	8.340000	84.00	83.229947
15	8.990000	85.46	85.784689
16	8.780000	84.35	84.959311
17	5.640000	73.19	72.617940
18	7.190000	78.21	78.710018
19	7.310000	77.59	79.181662
20	8.610000	83.87	84.291147
21	8.950000	85.15	85.627474

	study_hours	student_marks_original	student_marks_predicted
22	6.180000	72.96	74.740341
23	7.860000	80.72	81.343368
24	5.490000	73.61	72.028384
25	7.620000	79.53	80.400078
26	7.260000	78.17	78.985144
27	8.050000	79.63	82.090138
28	6.995949	76.83	77.947324
29	8.090000	82.38	82.247353
30	6.360000	76.04	75.447808
31	8.690000	85.48	84.605577
32	5.390000	71.87	71.635346
33	6.370000	75.04	75.487112
34	5.050000	70.67	70.299020
35	7.260000	79.87	78.985144
36	6.330000	74.49	75.329897
37	8.670000	84.19	84.526970
38	6.010000	75.36	74.072178
39	5.340000	72.10	71.438828

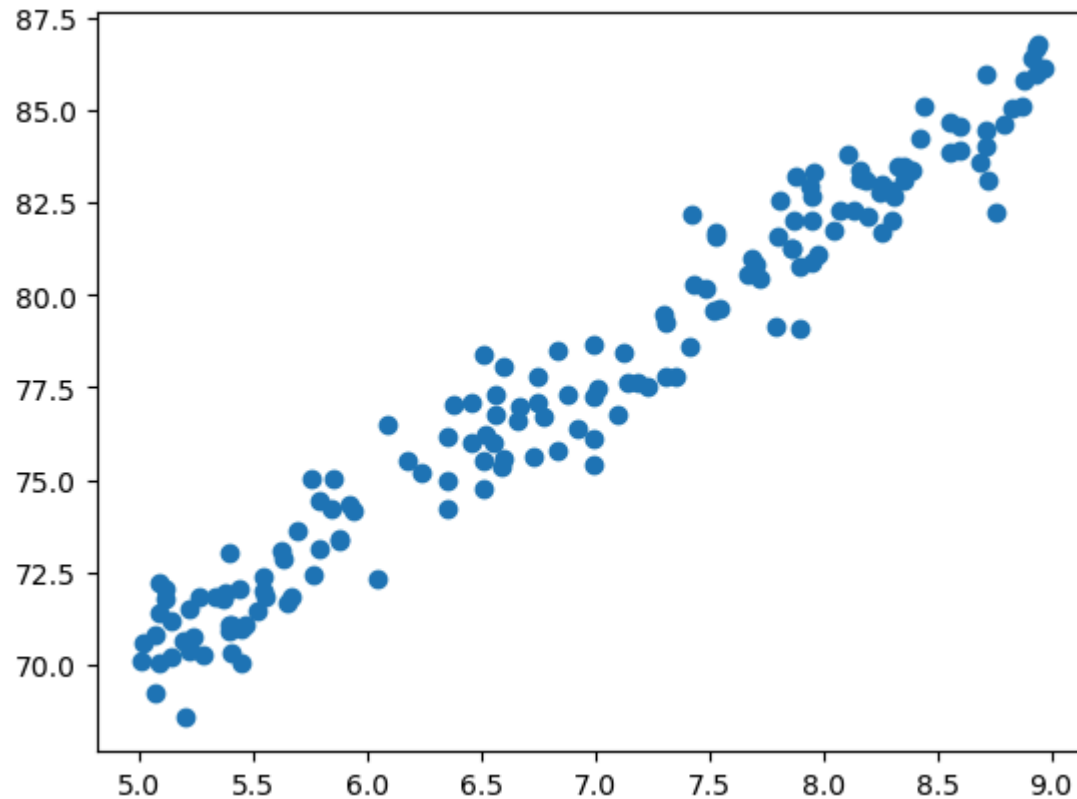
Fine-tune your model

```
In [52]: lr.score(x_test,y_test)
```

```
Out[52]: 0.9521841793508594
```

```
In [53]: plt.scatter(x_train,y_train)
```

```
Out[53]: <matplotlib.collections.PathCollection at 0x18b02583390>
```



```
plt.scatter(x_test, y_test) plt.plot(x_train, lr.predict(x_train), color = "r")
```

Present your solution

```
In [56]: import joblib  
joblib.dump(lr, "student_mark_predictor.pkl")
```

```
Out[56]: ['student_mark_predictor.pkl']
```

```
In [59]: pwd
```

```
Out[59]: 'C:\\Users\\samua'
```

```
In [57]: model = joblib.load("student_mark_predictor.pkl") # pickle
```

```
In [60]: pwd
```

```
Out[60]: 'C:\\Users\\samua'
```

```
In [58]: model.predict([[5]])[0][0]
```

```
C:\Users\samua\anaconda3\Lib\site-packages\sklearn\utils\validation.py:2739: UserWarning: X does not have valid feature names, but LinearRegression was fitted with feature names
  warnings.warn(
```

```
Out[58]: np.float64(70.10250100162845)
```

Launch, monitor, and maintain your system

```
In [61]: __name__
```

```
Out[61]: '__main__'
```

```
In [ ]:
```