

In [2]: `import pandas as pd
import numpy as np`

In [5]: `df = pd.read_csv(r"C:\Users\dheeman\Downloads\student_info.csv")`

In [6]: `df.head()`

Out[6]:

	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 [8]: `df.tail()`

Out[8]:

	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

In [11]: `df.shape`

Out[11]: (200, 2)

In [13]: `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: 3.2 KB
```

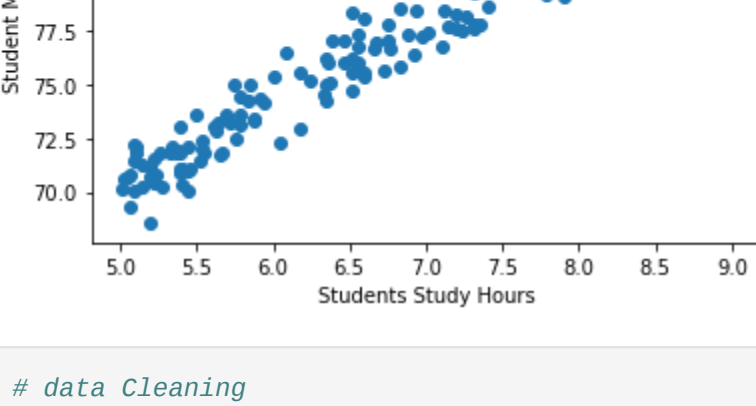
In [14]: `df.describe()`

Out[14]:

	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 [15]: `import matplotlib.pyplot as plt`

In [17]: `plt.scatter(x =df.study_hours, y = df. student_marks)
plt.xlabel("Students Study Hours")
plt.ylabel("Student Mark")
plt.title("scatter plot if student study hours vs Students mark")
plt.show()`



In [18]: `# data Cleaning`

In [19]: `df.isnull().sum()`

study\_hours 5  
student\_marks 0  
dtype: int64

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

study\_hours 6.995949  
student\_marks 77.933750  
dtype: float64

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

In [22]: `df2.isnull().sum()`

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 [25]: `df2.head()`

Out[25]:

	study_hours	student_marks
0	6.830000	78.50
1	6.560000	76.74
2	6.995949	78.68
3	5.670000	71.82
4	8.670000	84.19

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]: `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=51)
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)

## Select a model and train it

In [29]: `# Y = m * X + c
from sklearn.linear_model import LinearRegression
lr = LinearRegression()`

In [33]: `lr.fit (X_train,y_train)`

Out[33]: LinearRegression()

In [34]: `lr.coef_`

array([[3.93571802]])

In [35]: `lr.intercept_`

array([50.44735504])

In [36]: `m =3.93
c = 50.44
y = m * 4 + c
y`

66.16

In [37]: `lr.predict([[4]])[0].round(2)`

66.19

In [38]: `y_pred = lr.predict(X_test)
y_pred`

Out[38]: array([[83.11381458],
 [78.9025963 ],
 [84.57003024],
 [85.82946901],
 [84.72745896],
 [80.75238377],
 [72.84159055],
 [71.66087515],
 [73.23516235],
 [71.66087515],
 [73.47130543],
 [76.38373677],
 [73.23516235],
 [73.58937697],
 [82.95638585],
 [70.40144538],
 [73.23516235],
 [78.74516758],
 [75.55723598],
 [82.68088559],
 [76.65923703],
 [70.48015974],
 [74.77009238],
 [77.98143645],
 [85.59331693],
 [82.56281405],
 [76.42309395],
 [85.0423164 ],
 [78.39095296],
 [81.38209865],
 [81.73631327],
 [83.15317176],
 [82.20859943],
 [81.10659839],
 [73.58937697],
 [71.1492318 ],
 [71.89701823],
 [81.53952737],
 [72.60544747],
 [71.93637541]])

In [41]: `pd.DataFrame(np.c_[X_test,y_test, y_pred],columns = ["study_hours", "student_mark_orignal", "student_marks_predicte`

Out[41]:

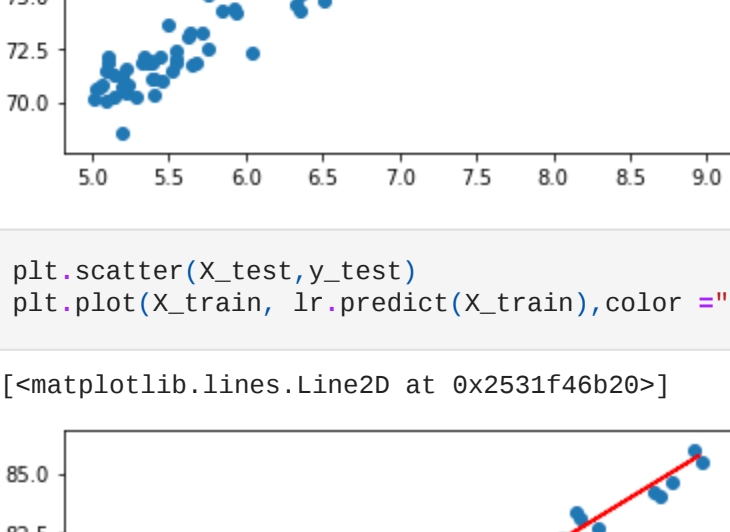
	study_hours	student_mark_orignal	student_marks_predicted
0	8.300000	82.02	83.113815
1	7.230000	77.55	78.902596
2	8.670000	84.19	84.570030
3	8.990000	85.46	85.829460
4	8.710000	84.03	84.727459
5	7.700000	80.81	80.752384
6	5.690000	73.61	72.841591
7	5.390000	70.90	71.660875
8	5.790000	73.14	73.235162
9	5.390000	73.02	71.660875
10	5.850000	75.02	73.471305
11	6.590000	75.37	76.383737
12	5.790000	74.44	73.235162
13	5.880000	73.40	73.589377
14	8.260000	81.70	82.956386
15	5.070000	69.27	70.401445
16	5.790000	73.64	73.235162
17	7.190000	77.63	78.745168
18	6.380000	77.01	75.557236
19	8.190000	83.08	82.680886
20	6.660000	76.63	76.659237
21	5.090000	72.22	70.480160
22	6.180000	72.96	74.770092
23	6.995949	76.14	77.981436
24	8.930000	85.96	85.593317
25	8.160000	83.36	82.562814
26	6.600000	78.05	76.423094
27	8.790000	84.60	85.042316
28	7.100000	76.76	78.390953
29	7.860000	81.24	81.382099
30	7.950000	80.86	81.736313
31	8.310000	82.69	83.153172
32	8.070000	82.30	82.208599
33	7.790000	79.17	81.106598
34	5.880000	73.34	73.589377
35	5.260000	71.86	71.149232
36	5.450000	70.06	71.897018
37	7.900000	80.76	81.539527
38	5.630000	72.87	72.605447
39	5.460000	71.10	71.936375

In [43]: `lr.score(X_test,y_test)`

Out[43]: 0.9514124242154464

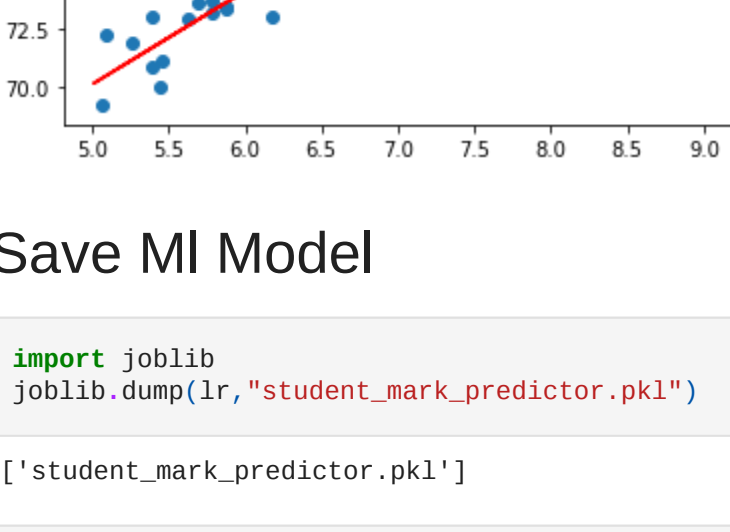
In [46]: `plt.scatter(X_train,y_train)`

Out[46]: <matplotlib.collections.PathCollection at 0x25398e4730>



In [48]: `plt.scatter(X_test,y_test)
plt.plot(X_train,lr.predict(X_train),color = "r")`

Out[48]: <matplotlib.lines.Line2D at 0x2531f46b20>



## Save ML Model

In [51]: `import joblib
joblib.dump(lr,"student_mark_predictor.pkl")`

Out[51]: ['student\_mark\_predictor.pkl']

In [52]: `model = joblib.load("student_mark_predictor.pkl")`

In [53]: `model.predict([[5]])[0][0]`

Out[53]: 70.12594512018407

In [ ]: