

17-studentsmarksprediction

September 21, 2025

1 STUDENTS MARKS PREDICTION

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

```
[2]: df = pd.read_csv(r"D:\Naresh It Classes\4. September\17- Production, Into to AI\Students Marks Prediction\student_info.csv")
```

```
[3]: df.head()
```

```
[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

```
[4]: 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
```

```
[6]: df.shape
```

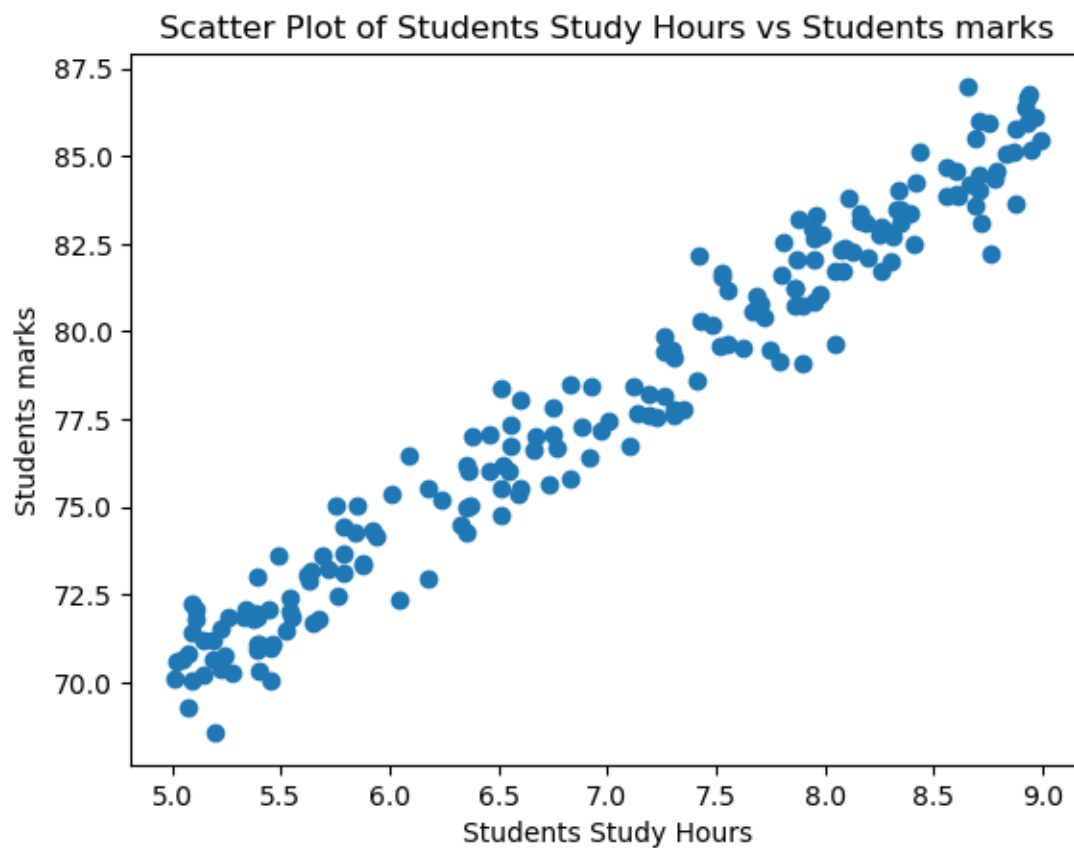
```
[6]: (200, 2)
```

```
[7]: df.describe()
```

```
[7]:
```

	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

```
[8]: 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()
```



1.1 Data Cleaning

```
[9]: df.isna()
```

```
[9]:      study_hours  student_marks
0           False           False
1           False           False
2            True           False
3           False           False
4           False           False
..           ...             ...
195          False           False
196          False           False
197          False           False
198          False           False
199          False           False
```

[200 rows x 2 columns]

```
[10]: df.isna().sum()
```

```
[10]: study_hours      5
      student_marks    0
      dtype: int64
```

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

```
[12]: df2.isna().sum()
```

```
[12]: study_hours      0
      student_marks    0
      dtype: int64
```

1.2 Split data

```
[13]: 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)
```

```
[14]: 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)
```

1.3 Choosing Linear Regression model and training it

```
[17]: from sklearn.linear_model import LinearRegression
lr = LinearRegression()
```

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

```
[18]: LinearRegression()
```

```
[24]: m = lr.coef_
m
```

```
[24]: array([[3.93037294]])
```

```
[25]: c = lr.intercept_
c
```

```
[25]: array([50.45063632])
```

```
[27]: y = m*4 + c #for 4 hours of study
y
```

```
[27]: array([[66.17212807]])
```

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

```
D:\Program Files\Lib\site-packages\sklearn\base.py:493: UserWarning: X does not
have valid feature names, but LinearRegression was fitted with feature names
warnings.warn(
```

```
[28]: 66.17
```

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

```
[29]: 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 ]])
```

```
[30]: pd.DataFrame(np.c_[X_test, y_test, y_pred], columns = ["study_hours",
↪ "student_marks_original", "student_marks_predicted"])
```

```
[30]:
```

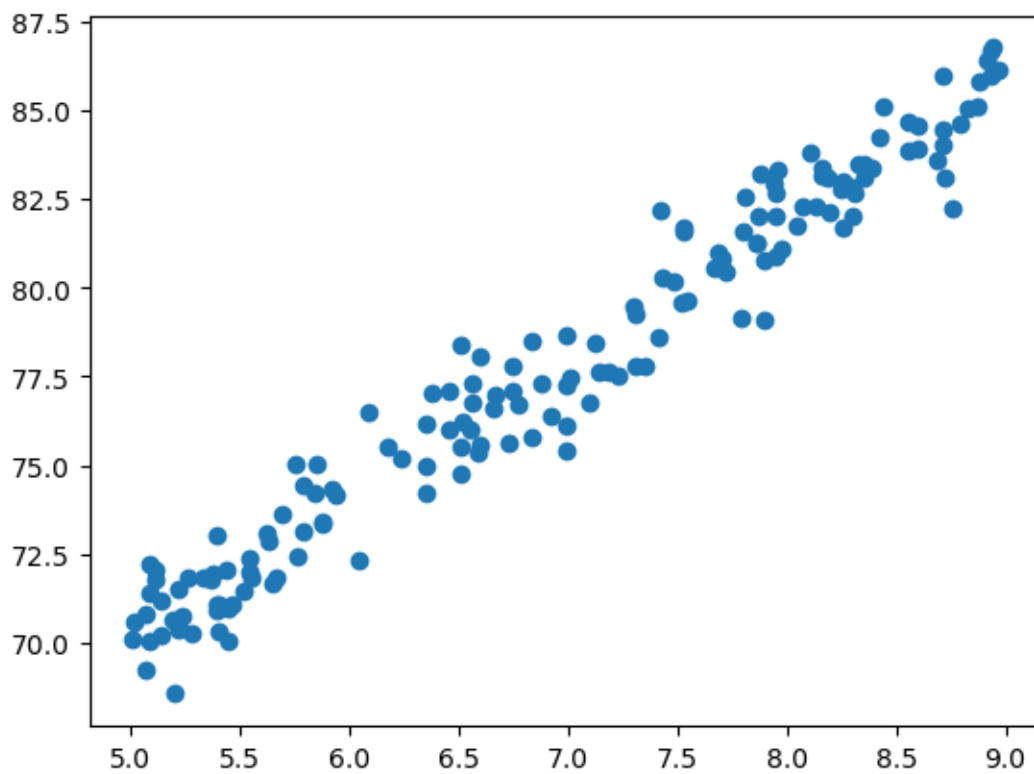
	study_hours	student_marks_original	student_marks_predicted
0	8.410000	82.50	83.505073
1	5.190000	71.18	70.849272
2	5.720000	73.25	72.932370
3	8.880000	83.64	85.352348
4	5.790000	73.64	73.207496

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
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

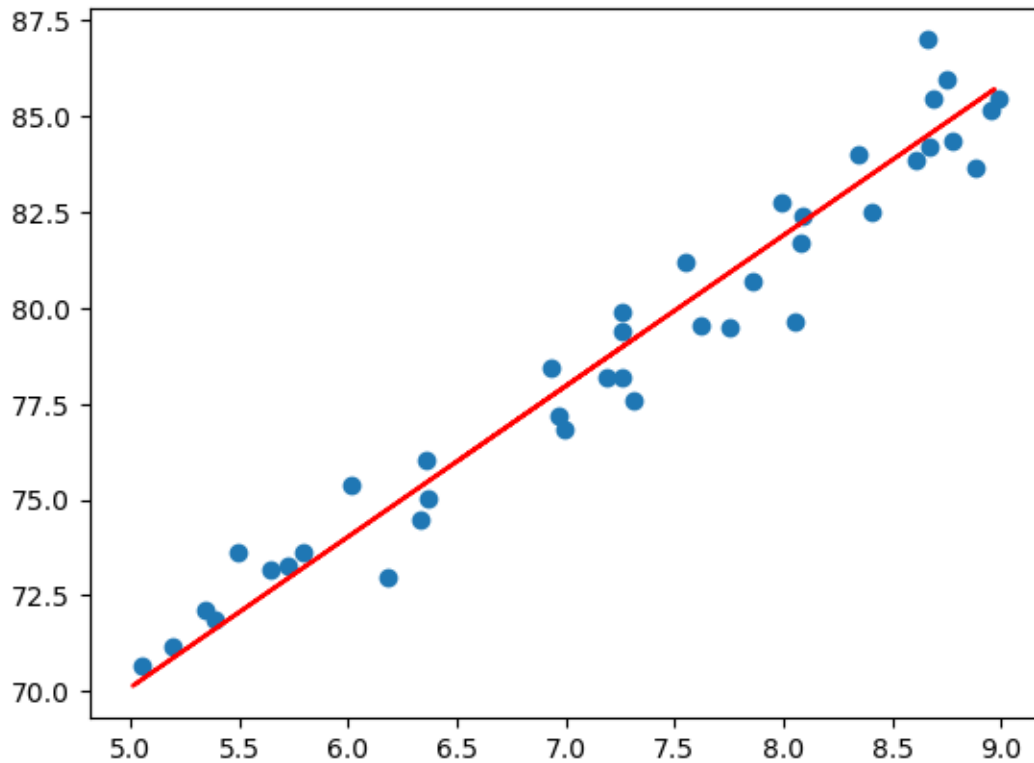
```
[31]: lr.score(X_test,y_test)    #R2
```

```
[31]: 0.9521841793508595
```

```
[35]: plt.scatter(X_train,y_train)
plt.show()
```



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



1.4 Save the model as a Pickle File

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

```
[36]: ['student_mark_predictor.pkl']
```

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

```
[38]: model.predict([[5]])[0][0]
```

D:\Program Files\Lib\site-packages\sklearn\base.py:493: UserWarning: X does not have valid feature names, but LinearRegression was fitted with feature names
warnings.warn(

```
[38]: 70.10250100162847
```

```
[ ]:
```