

Problem Statement

Linear Regression

Import Libraries

In [1]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

In [2]:

```
a=pd.read_csv("student.csv")
a
```

Out[2]:

	Student_ID	Test_1	Test_2	Test_3	Test_4	Test_5	Test_6	Test_7	Test_8	Test_9	Test_10	Test_11
0	22000	78	87	91	91	88	98	94	100	100	100	100
1	22001	79	71	81	72	73	68	59	69	59	60	60
2	22002	66	65	70	74	78	86	87	96	88	82	90
3	22003	60	58	54	61	54	57	64	62	72	63	70
4	22004	99	95	96	93	97	89	92	98	91	98	99
5	22005	41	36	35	28	35	36	27	26	19	22	20
6	22006	47	50	47	57	62	64	71	75	85	87	88
7	22007	84	74	70	68	58	59	56	56	64	70	60
8	22008	74	64	58	57	53	51	47	45	42	43	30
9	22009	87	81	73	74	71	63	53	45	39	43	40
10	22010	40	34	37	33	31	35	39	38	40	48	40
11	22011	91	84	78	74	76	80	80	73	75	71	70
12	22012	81	83	93	88	89	90	99	99	95	85	70
13	22013	52	50	42	38	33	30	28	22	12	20	10
14	22014	63	67	65	74	80	86	95	96	92	83	70
15	22015	76	82	88	94	85	76	70	60	50	58	40
16	22016	83	78	71	71	77	72	66	75	66	61	60
17	22017	55	45	43	38	43	35	44	37	45	37	40
18	22018	71	67	76	74	64	61	57	64	61	51	50
19	22019	62	61	53	49	54	59	68	74	65	55	60
20	22020	44	38	36	34	26	34	39	44	36	45	30
21	22021	50	56	53	46	41	38	47	39	44	36	40
22	22022	57	48	40	45	43	36	26	19	9	12	20

	Student_ID	Test_1	Test_2	Test_3	Test_4	Test_5	Test_6	Test_7	Test_8	Test_9	Test_10	Test_11
23	22023	59	56	52	44	50	40	45	46	54	57	51
24	22024	84	92	89	80	90	80	84	74	68	73	81
25	22025	74	80	86	87	90	100	95	87	85	79	83
26	22026	92	84	74	83	93	83	75	82	81	73	70
27	22027	63	70	74	65	64	55	61	58	48	46	41
28	22028	78	77	69	76	78	74	67	69	78	68	61
29	22029	55	58	59	67	71	62	53	61	67	76	72
30	22030	54	54	48	38	35	45	46	47	41	37	31
31	22031	84	93	97	89	86	95	100	100	100	99	100
32	22032	95	100	94	100	98	99	100	90	80	84	75
33	22033	64	61	63	73	63	68	64	58	50	51	51
34	22034	76	79	73	77	83	86	95	89	90	95	100
35	22035	78	71	61	55	54	48	41	32	41	40	41
36	22036	95	89	91	84	89	94	85	91	100	100	100
37	22037	99	89	79	87	87	81	82	74	64	54	51
38	22038	82	83	85	86	89	80	88	95	87	93	90
39	22039	65	56	64	62	58	51	61	68	70	70	61
40	22040	100	93	92	86	84	76	82	74	79	72	71
41	22041	78	72	73	79	81	73	71	77	83	92	91
42	22042	98	100	100	93	94	92	100	100	98	94	91
43	22043	58	62	67	77	71	63	64	73	83	76	80
44	22044	96	92	94	100	99	95	98	92	84	84	81
45	22045	86	87	85	84	85	91	86	82	85	87	81
46	22046	48	55	46	40	34	29	37	34	39	41	31
47	22047	56	52	54	47	40	35	43	44	40	39	41
48	22048	42	44	46	53	62	59	57	53	43	35	31
49	22049	64	54	49	59	54	55	57	59	63	73	71
50	22050	50	44	37	29	37	46	53	57	55	61	61
51	22051	70	60	70	62	67	67	68	67	72	69	61
52	22052	63	73	70	63	60	67	61	59	52	58	51
53	22053	92	100	100	100	100	100	92	87	94	100	91
54	22054	64	55	54	61	63	57	47	37	44	48	51
55	22055	60	66	68	58	49	47	39	29	39	44	31

To display top 10 rows

In [3]:

```
c=a.head(15)
c
```

Out[3]:

	Student_ID	Test_1	Test_2	Test_3	Test_4	Test_5	Test_6	Test_7	Test_8	Test_9	Test_10	Test_11
0	22000	78	87	91	91	88	98	94	100	100	100	100
1	22001	79	71	81	72	73	68	59	69	59	60	60
2	22002	66	65	70	74	78	86	87	96	88	82	90
3	22003	60	58	54	61	54	57	64	62	72	63	70
4	22004	99	95	96	93	97	89	92	98	91	98	90
5	22005	41	36	35	28	35	36	27	26	19	22	20
6	22006	47	50	47	57	62	64	71	75	85	87	80
7	22007	84	74	70	68	58	59	56	56	64	70	60
8	22008	74	64	58	57	53	51	47	45	42	43	30
9	22009	87	81	73	74	71	63	53	45	39	43	40
10	22010	40	34	37	33	31	35	39	38	40	48	40
11	22011	91	84	78	74	76	80	80	73	75	71	70
12	22012	81	83	93	88	89	90	99	99	95	85	70
13	22013	52	50	42	38	33	30	28	22	12	20	10
14	22014	63	67	65	74	80	86	95	96	92	83	70

To find Missing values

In [4]:

```
c.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 15 entries, 0 to 14
Data columns (total 13 columns):
#   Column      Non-Null Count  Dtype  
---  -
0   Student_ID  15 non-null    int64  
1   Test_1      15 non-null    int64  
2   Test_2      15 non-null    int64  
3   Test_3      15 non-null    int64  
4   Test_4      15 non-null    int64  
5   Test_5      15 non-null    int64  
6   Test_6      15 non-null    int64  
7   Test_7      15 non-null    int64  
8   Test_8      15 non-null    int64  
9   Test_9      15 non-null    int64  
10  Test_10     15 non-null    int64  
11  Test_11     15 non-null    int64  
12  Test_12     15 non-null    int64  
dtypes: int64(13)
memory usage: 1.6 KB
```

To display summary of statistics

In [5]:

a.describe()

Out[5]:

	Student_ID	Test_1	Test_2	Test_3	Test_4	Test_5	Test_6	Te
count	56.000000	56.000000	56.000000	56.000000	56.000000	56.000000	56.000000	56.000
mean	22027.500000	70.750000	69.196429	68.089286	67.446429	67.303571	66.000000	66.160
std	16.309506	17.009356	17.712266	18.838333	19.807179	20.746890	21.054043	21.427
min	22000.000000	40.000000	34.000000	35.000000	28.000000	26.000000	29.000000	26.000
25%	22013.750000	57.750000	55.750000	53.000000	54.500000	53.750000	50.250000	47.000
50%	22027.500000	70.500000	68.500000	70.000000	71.500000	69.000000	65.500000	64.000
75%	22041.250000	84.000000	83.250000	85.000000	84.000000	85.250000	83.750000	85.250
max	22055.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000

To display column heading

In [6]:

a.columns

Out[6]:

Index(['Student_ID', 'Test_1', 'Test_2', 'Test_3', 'Test_4', 'Test_5',
 'Test_6', 'Test_7', 'Test_8', 'Test_9', 'Test_10', 'Test_11',
 'Test_12'],
 dtype='object')

Pairplot

In [7]:

s=a.dropna(axis=1)
s

Out[7]:

	Student_ID	Test_1	Test_2	Test_3	Test_4	Test_5	Test_6	Test_7	Test_8	Test_9	Test_10	Test_11
0	22000	78	87	91	91	88	98	94	100	100	100	100
1	22001	79	71	81	72	73	68	59	69	59	60	60
2	22002	66	65	70	74	78	86	87	96	88	82	90
3	22003	60	58	54	61	54	57	64	62	72	63	70
4	22004	99	95	96	93	97	89	92	98	91	98	90
5	22005	41	36	35	28	35	36	27	26	19	22	20
6	22006	47	50	47	57	62	64	71	75	85	87	80
7	22007	84	74	70	68	58	59	56	56	64	70	60
8	22008	74	64	58	57	53	51	47	45	42	43	30
9	22009	87	81	73	74	71	63	53	45	39	43	40
10	22010	40	34	37	33	31	35	39	38	40	48	40
11	22011	91	84	78	74	76	80	80	73	75	71	70
12	22012	81	83	93	88	89	90	99	99	95	85	70

	Student_ID	Test_1	Test_2	Test_3	Test_4	Test_5	Test_6	Test_7	Test_8	Test_9	Test_10	Test_11
13	22013	52	50	42	38	33	30	28	22	12	20	19
14	22014	63	67	65	74	80	86	95	96	92	83	79
15	22015	76	82	88	94	85	76	70	60	50	58	49
16	22016	83	78	71	71	77	72	66	75	66	61	60
17	22017	55	45	43	38	43	35	44	37	45	37	40
18	22018	71	67	76	74	64	61	57	64	61	51	50
19	22019	62	61	53	49	54	59	68	74	65	55	60
20	22020	44	38	36	34	26	34	39	44	36	45	30
21	22021	50	56	53	46	41	38	47	39	44	36	40
22	22022	57	48	40	45	43	36	26	19	9	12	20
23	22023	59	56	52	44	50	40	45	46	54	57	50
24	22024	84	92	89	80	90	80	84	74	68	73	80
25	22025	74	80	86	87	90	100	95	87	85	79	80
26	22026	92	84	74	83	93	83	75	82	81	73	70
27	22027	63	70	74	65	64	55	61	58	48	46	40
28	22028	78	77	69	76	78	74	67	69	78	68	60
29	22029	55	58	59	67	71	62	53	61	67	76	70
30	22030	54	54	48	38	35	45	46	47	41	37	30
31	22031	84	93	97	89	86	95	100	100	100	99	100
32	22032	95	100	94	100	98	99	100	90	80	84	70
33	22033	64	61	63	73	63	68	64	58	50	51	50
34	22034	76	79	73	77	83	86	95	89	90	95	100
35	22035	78	71	61	55	54	48	41	32	41	40	40
36	22036	95	89	91	84	89	94	85	91	100	100	100
37	22037	99	89	79	87	87	81	82	74	64	54	50
38	22038	82	83	85	86	89	80	88	95	87	93	90
39	22039	65	56	64	62	58	51	61	68	70	70	60
40	22040	100	93	92	86	84	76	82	74	79	72	70
41	22041	78	72	73	79	81	73	71	77	83	92	90
42	22042	98	100	100	93	94	92	100	100	98	94	90
43	22043	58	62	67	77	71	63	64	73	83	76	80
44	22044	96	92	94	100	99	95	98	92	84	84	80
45	22045	86	87	85	84	85	91	86	82	85	87	80
46	22046	48	55	46	40	34	29	37	34	39	41	30
47	22047	56	52	54	47	40	35	43	44	40	39	40
48	22048	42	44	46	53	62	59	57	53	43	35	30

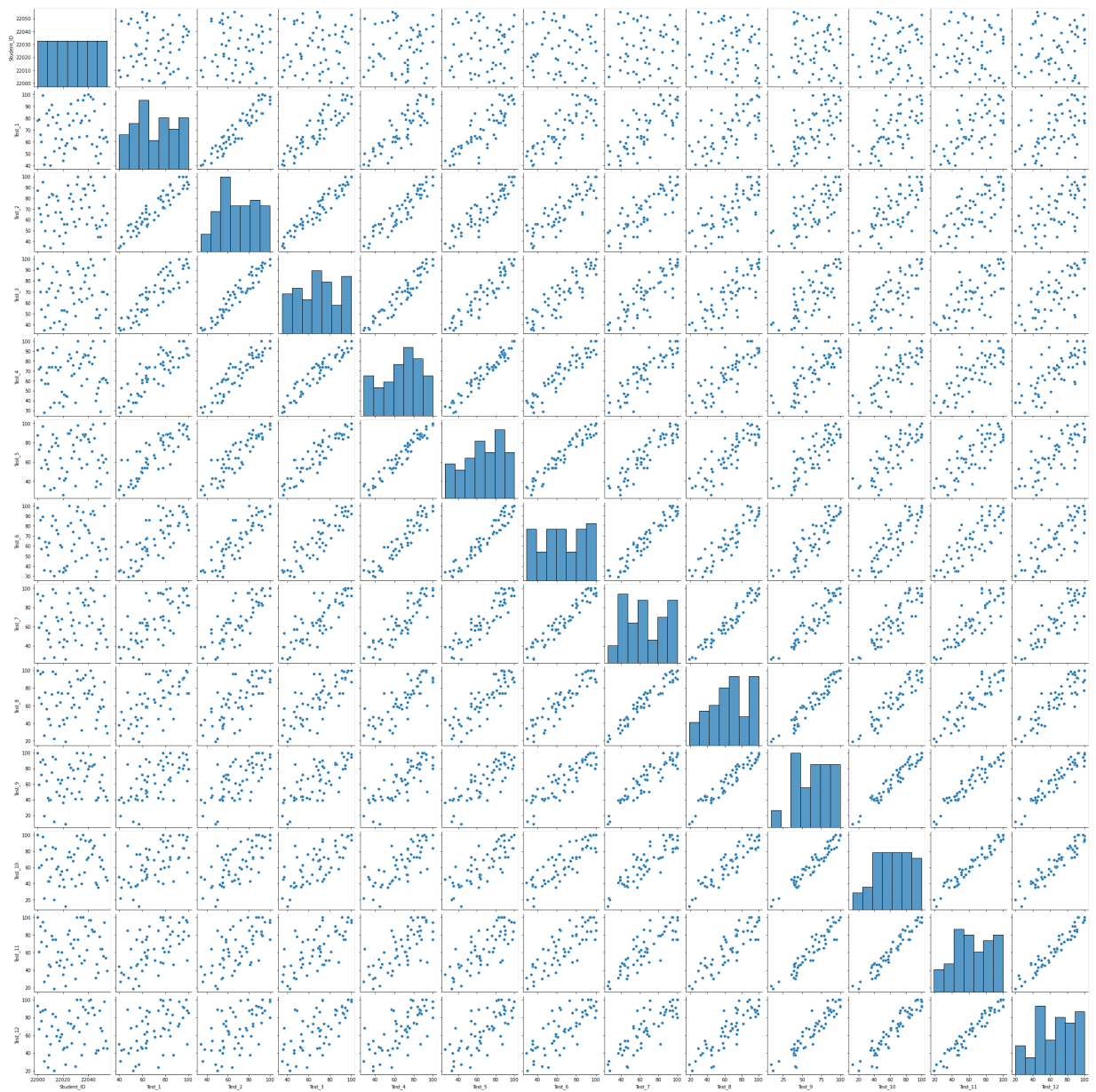
	Student_ID	Test_1	Test_2	Test_3	Test_4	Test_5	Test_6	Test_7	Test_8	Test_9	Test_10	Test_11
49	22049	64	54	49	59	54	55	57	59	63	73	71
50	22050	50	44	37	29	37	46	53	57	55	61	60
51	22051	70	60	70	62	67	67	68	67	72	69	60
52	22052	63	73	70	63	60	67	61	59	52	58	51
53	22053	92	100	100	100	100	100	92	87	94	100	90
54	22054	64	55	54	61	63	57	47	37	44	48	50
55	22055	60	66	68	58	49	47	39	29	39	44	31

In [8]: `s.columns`

Out[8]: Index(['Student_ID', 'Test_1', 'Test_2', 'Test_3', 'Test_4', 'Test_5',
'Test_6', 'Test_7', 'Test_8', 'Test_9', 'Test_10', 'Test_11',
'Test_12'],
dtype='object')

In [10]: `sns.pairplot(a)`

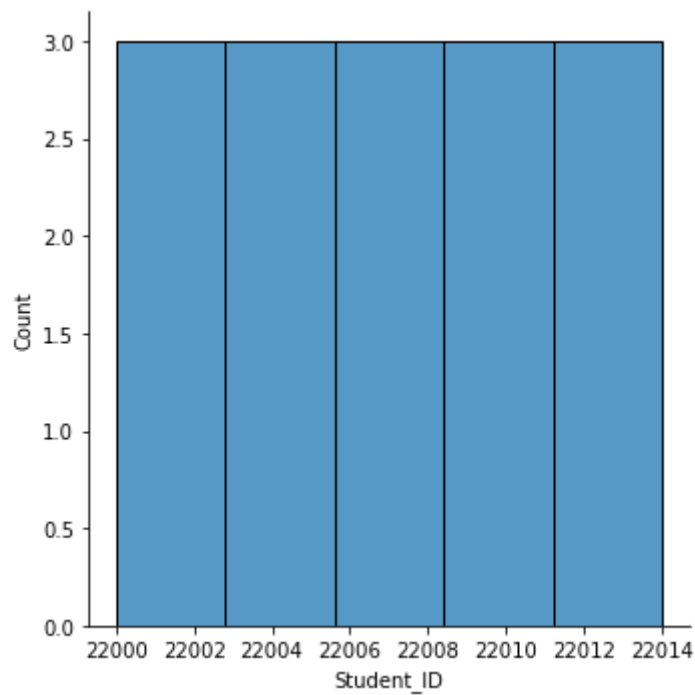
Out[10]: <seaborn.axisgrid.PairGrid at 0x2b08c2dba90>



Distribution Plot

```
In [11]: sns.displot(c['Student_ID'])
```

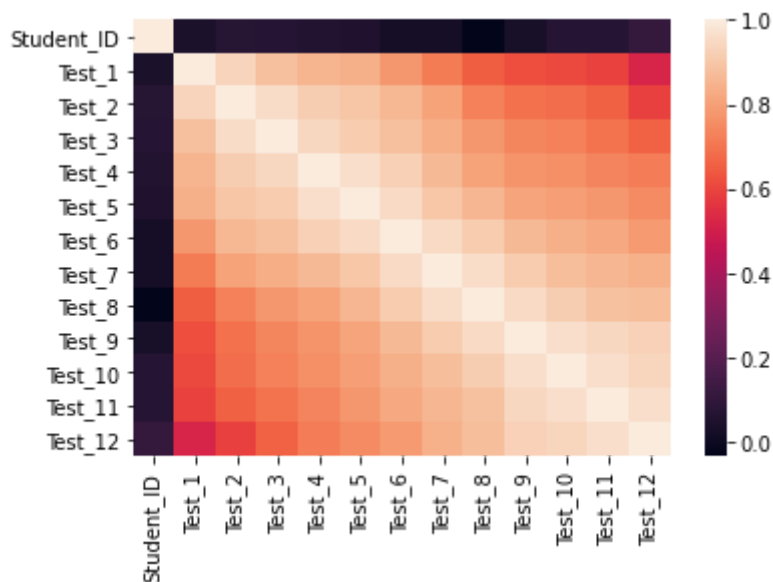
```
Out[11]: <seaborn.axisgrid.FacetGrid at 0x2b08c3a2b50>
```



Correlation

```
In [12]: b=a[['Student_ID', 'Test_1', 'Test_2', 'Test_3', 'Test_4', 'Test_5',
              'Test_6', 'Test_7', 'Test_8', 'Test_9', 'Test_10', 'Test_11',
              'Test_12']]
          sns.heatmap(b.corr())
```

Out[12]: <AxesSubplot:>



Train the model - Model Building

```
In [13]: g=c[['Student_ID', 'Test_1', 'Test_2', 'Test_3', 'Test_4', 'Test_5',
              'Test_6', 'Test_7', 'Test_8', 'Test_9', 'Test_10', 'Test_11']]
          h=c['Test_12']
```


To split dataset into training end test

```
In [14]: from sklearn.model_selection import train_test_split
g_train,g_test,h_train,h_test=train_test_split(g,h,test_size=0.6)
```

To run the model

```
In [15]: from sklearn.linear_model import LinearRegression
```

```
In [16]: lr=LinearRegression()
lr.fit(g_train,h_train)
```

Out[16]: LinearRegression()

```
In [17]: print(lr.intercept_)
```

6794.729390518639

Coeffecient

```
In [18]: coeff=pd.DataFrame(lr.coef_,g.columns,columns=['Co-effecient'])
coeff
```

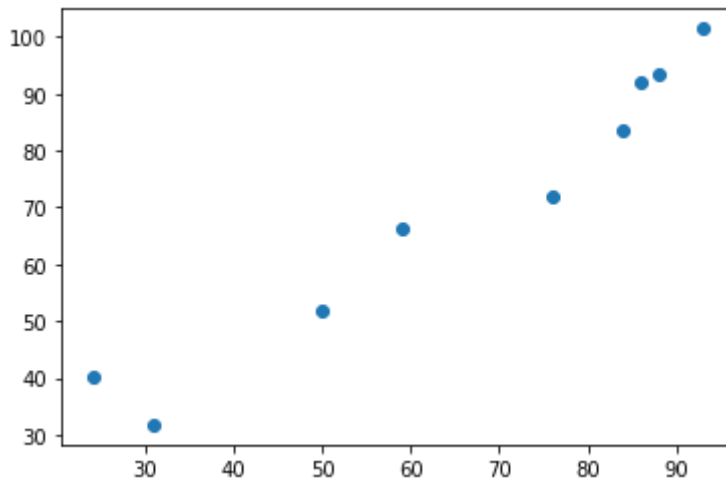
Out[18]:

	Co-effecient
Student_ID	-0.308044
Test_1	-0.059771
Test_2	-0.220402
Test_3	0.268847
Test_4	-0.122618
Test_5	-0.065898
Test_6	-0.065108
Test_7	-0.171812
Test_8	0.300793
Test_9	0.255367
Test_10	0.290472
Test_11	0.378750

Best Fit line

```
In [19]: prediction=lr.predict(g_test)
plt.scatter(h_test,prediction)
```

Out[19]: <matplotlib.collections.PathCollection at 0x2b0975745e0>



To find score

```
In [20]: print(lr.score(g_test,h_test))
```

0.9114400373695808

Import Lasso and ridge

```
In [21]: from sklearn.linear_model import Ridge,Lasso
```

Ridge

```
In [22]: ri=Ridge(alpha=5)
ri.fit(g_train,h_train)
```

Out[22]: Ridge(alpha=5)

```
In [23]: ri.score(g_test,h_test)
```

Out[23]: 0.912333818759164

```
In [24]: ri.score(g_train,h_train)
```

Out[24]: 0.9999883017124321

Lasso

```
In [25]: l=Lasso(alpha=6)
l.fit(g_train,h_train)
```

Out[25]: Lasso(alpha=6)

In [26]: `l.score(g_test,h_test)`

Out[26]: 0.9239353234554334

In [27]: `ri.score(g_train,h_train)`

Out[27]: 0.9999883017124321

In []: