

Problem Statement

Problem Statement

```
import pandas as pd
```

```
df.drop(columns=['Serial No.'], inplace=True)
print(df.shape)
```

```
print(df.duplicated().sum())
```

DOA	0
OGPA	0

```
In [5]: df.describe(include='all')
```

1118A	390,000,000	1,200,000,000	1118Y	0,000,000	0,000,000	0,000,000	1118Y	0,000,000
-------	-------------	---------------	-------	-----------	-----------	-----------	-------	-----------

```
plt.title(f'Distribution of {col}')
plt.show()
```

Age Group	Number of People
0-10	100
11-20	90
21-30	80
31-40	70
41-50	60
51-60	50
61-70	40
71-80	30
81-90	20
91-100	10

Category	Count
no	1
yes	1

Category	Number of people
Do not use the Internet	10

Distribution of LOR	
1	0.0000
2	0.0000
3	0.0000
4	0.0000
5	0.0000
6	0.0000
7	0.0000
8	0.0000
9	0.0000
10	0.0000
11	0.0000
12	0.0000
13	0.0000
14	0.0000
15	0.0000
16	0.0000
17	0.0000
18	0.0000
19	0.0000
20	0.0000
21	0.0000
22	0.0000
23	0.0000
24	0.0000
25	0.0000
26	0.0000
27	0.0000
28	0.0000
29	0.0000
30	0.0000
31	0.0000
32	0.0000
33	0.0000
34	0.0000
35	0.0000
36	0.0000
37	0.0000
38	0.0000
39	0.0000
40	0.0000
41	0.0000
42	0.0000
43	0.0000
44	0.0000
45	0.0000
46	0.0000
47	0.0000
48	0.0000
49	0.0000
50	0.0000
51	0.0000
52	0.0000
53	0.0000
54	0.0000
55	0.0000
56	0.0000
57	0.0000
58	0.0000
59	0.0000
60	0.0000
61	0.0000
62	0.0000
63	0.0000
64	0.0000
65	0.0000
66	0.0000
67	0.0000
68	0.0000
69	0.0000
70	0.0000
71	0.0000
72	0.0000
73	0.0000
74	0.0000
75	0.0000
76	0.0000
77	0.0000
78	0.0000
79	0.0000
80	0.0000
81	0.0000
82	0.0000
83	0.0000
84	0.0000
85	0.0000
86	0.0000
87	0.0000
88	0.0000
89	0.0000
90	0.0000
91	0.0000
92	0.0000
93	0.0000
94	0.0000
95	0.0000
96	0.0000
97	0.0000
98	0.0000
99	0.0000
100	0.0000

Year	Number of people
2018	60
2019	70

Chance of Admit

Country	Number of people
United States	140
Other countries	100

Year	Percentage
1990	38
1992	42
1994	45
1996	48
1998	50
2000	52
2002	50
2004	52

[illegible][illegible]

```
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
```

LOR	0.52	0.54	0.66	1	0.64	0.65
-----	------	------	------	---	------	------

Y

Scatter plot showing the normalized change in relative abundance of the 16S rRNA gene over time. The y-axis is labeled 'change' and ranges from 0.0 to 0.6. The x-axis is labeled 'time' and ranges from 0 to 10. The data points are blue dots, showing a general increase in change over time, with a peak around time 5 and then a slight decrease.

```
X_sm = add_constant(X)
model = OLS(y, X_sm).fit()
```

	coef	std err	t	P> t
--	------	---------	---	------

[illegible]

TOEFL Score	TOEFL Score	0.002778
University Rating	University Rating	0.003941

3	SOP	35.265006
4	Ende	30.033876

Circumstance	Percentage (%)
Self-defense	85
To protect others	75
To protect property	65
To protect the community	55
To protect the environment	45

QQ Plot

```
y_pred_test = model_train.predict(X_test)
```

Testing:
MAE: 0.0427, RMSE: 0.0609, R2: 0.8188, Adj R2: 0.8029

```
Lasso Coefficients: [0.02009774 0.0139114 0.0043499 0. 0
```

