

Jamboree Education

1 Jamboree has assisted numerous students in achieving admission to prestigious international colleges. Whether it's GMAT, GRE, or SAT, their distinctive problem-solving approaches guarantee optimal scores with minimal effort. They have introduced a new feature on their website allowing students to assess their likelihood of acceptance into Ivy League colleges. This feature provides an estimation of graduate admission probabilities tailored to an Indian perspective.

1 This case study involves utilizing the Jamboree Education dataset to analyze and forecast students' admission probabilities to leading global universities. The objective is to examine and predict the likelihood of students gaining acceptance into top-tier academic institutions worldwide.

Importing Libraries

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib
5 import matplotlib.pyplot as plt
6 from sklearn.model_selection import train_test_split
7 from sklearn.preprocessing import MinMaxScaler
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.linear_model import LinearRegression
10 from sklearn.metrics import r2_score
11 import statsmodels.api as sm
12 from statsmodels.stats.outliers_influence import variance_inflation_factor
13 from sklearn.metrics import mean_squared_error
14 from sklearn.metrics import mean_absolute_error
15 import statsmodels.api as sm
16 from sklearn.preprocessing import PolynomialFeatures
17 from sklearn.linear_model import Lasso , Ridge
```

```
In [2]: 1 df = pd.read_csv('Jamboree_Admission.csv')
2 df
```

Out[2]:

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
0	1	337	118	4	4.5	4.5	9.65	1	0.92
1	2	324	107	4	4.0	4.5	8.87	1	0.76
2	3	316	104	3	3.0	3.5	8.00	1	0.72
3	4	322	110	3	3.5	2.5	8.67	1	0.80
4	5	314	103	2	2.0	3.0	8.21	0	0.65
...
495	496	332	108	5	4.5	4.0	9.02	1	0.87
496	497	337	117	5	5.0	5.0	9.87	1	0.96
497	498	330	120	5	4.5	5.0	9.56	1	0.93
498	499	312	103	4	4.0	5.0	8.43	0	0.73
499	500	327	113	4	4.5	4.5	9.04	0	0.84

500 rows × 9 columns

```
In [3]: 1 df.head()
```

Out[3]:

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
0	1	337	118	4	4.5	4.5	9.65	1	0.92
1	2	324	107	4	4.0	4.5	8.87	1	0.76
2	3	316	104	3	3.0	3.5	8.00	1	0.72
3	4	322	110	3	3.5	2.5	8.67	1	0.80
4	5	314	103	2	2.0	3.0	8.21	0	0.65

```
In [4]: 1 df.shape
```

Out[4]: (500, 9)

```
In [5]: 1 df.size
```

Out[5]: 4500

```
In [6]: 1 pd.DataFrame(df.dtypes)
```

Out[6]:

	0
Serial No.	int64
GRE Score	int64
TOEFL Score	int64
University Rating	int64
SOP	float64
LOR	float64
CGPA	float64
Research	int64
Chance of Admit	float64

```
In [7]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Serial No.       500 non-null    int64  
 1   GRE Score        500 non-null    int64  
 2   TOEFL Score      500 non-null    int64  
 3   University Rating 500 non-null    int64  
 4   SOP              500 non-null    float64 
 5   LOR              500 non-null    float64 
 6   CGPA             500 non-null    float64 
 7   Research          500 non-null    int64  
 8   Chance of Admit  500 non-null    float64 
dtypes: float64(4), int64(5)
memory usage: 35.3 KB
```

```
In [8]: 1 df.isna().sum()
```

```
Out[8]: Serial No.      0
GRE Score       0
TOEFL Score     0
University Rating 0
SOP             0
LOR             0
CGPA            0
Research         0
Chance of Admit  0
dtype: int64
```

1 This Data does not have any null values present, every feild contain some values, that will be better to Interpretate the segments without harm of the data.

In [9]: 1 df.rename(columns={'LOR':"LOR"}, inplace=True)

In [10]: 1 df.drop(columns={"Serial No."}, inplace=True)

In [11]: 1 df.duplicated().sum()

Out[11]: 0

1 In the given dataset, there are no duplicated items.

In [12]: 1 df.columns

Out[12]: Index(['GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR ', 'CGPA',
'Research', 'Chance of Admit '],
dtype='object')

In [13]: 1 df.describe()

Out[13]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Cl of.
count	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.
mean	316.472000	107.192000	3.114000	3.374000	3.48400	8.576440	0.560000	0.
std	11.295148	6.081868	1.143512	0.991004	0.92545	0.604813	0.496884	0.
min	290.000000	92.000000	1.000000	1.000000	1.00000	6.800000	0.000000	0.
25%	308.000000	103.000000	2.000000	2.500000	3.00000	8.127500	0.000000	0.
50%	317.000000	107.000000	3.000000	3.500000	3.50000	8.560000	1.000000	0.
75%	325.000000	112.000000	4.000000	4.000000	4.00000	9.040000	1.000000	0.
max	340.000000	120.000000	5.000000	5.000000	5.00000	9.920000	1.000000	0.

1 This is a comprehensive data summary incorporating optimal values for each feature within the dataset. It includes key statistical measures, such as quantiles ranging from 25% to 75%, to assess the potential impact of outliers and related scenarios.

Analysis Univariate / Bivariate

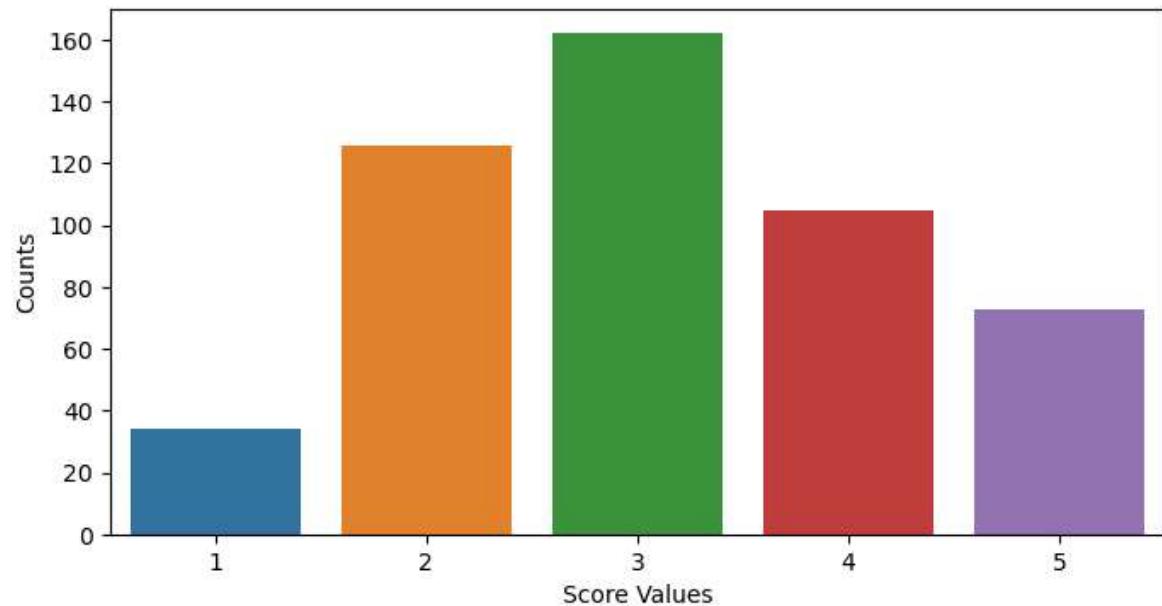
University_Rating

```
In [14]: 1 University_Rating = pd.DataFrame(df['University Rating'].value_counts(asc  
2 University_Rating
```

Out[14]:

index	University Rating
0	3
1	2
2	4
3	5
4	1

```
In [15]: 1 plt.figure(figsize=(8,4))  
2 sns.barplot(data=University_Rating, x='index', y='University Rating')  
3 plt.xlabel("Score Values")  
4 plt.ylabel("Counts")  
5 plt.show()
```



- According to the data, student ratings for the university, ranging from 1 to 5, were analyzed. The graph illustrates that the majority, approximately 162 out of 500 students, gave a rating of 3. Additionally, 73 students provided the highest rating of 5.

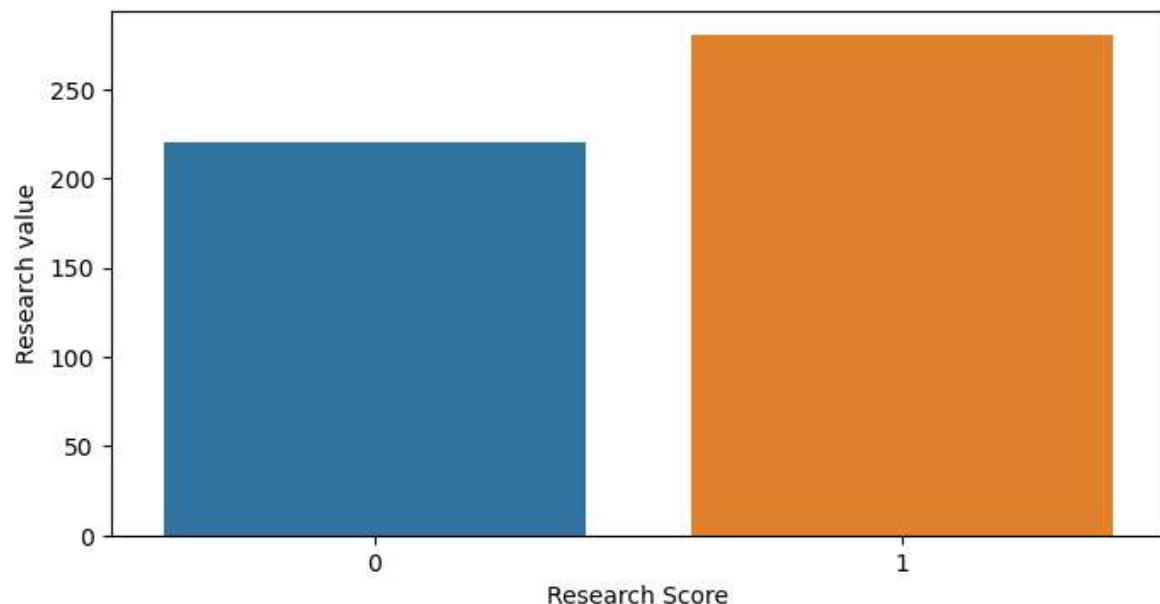
Research

```
In [16]: 1 Research = pd.DataFrame(df['Research'].value_counts(ascending=False)).reset_index()
          2 Research
```

Out[16]:

index	Research
0	1 280
1	0 220

```
In [17]: 1 plt.figure(figsize=(8,4))
          2 sns.barplot(data=Research, x='index', y='Research')
          3 plt.xlabel("Research Score")
          4 plt.ylabel("Research value")
          5 plt.show()
```



- This graph depicts the correlation between research experiences and students, highlighting the frequency of observations within this relationship.

SOP - Statement of Purpose

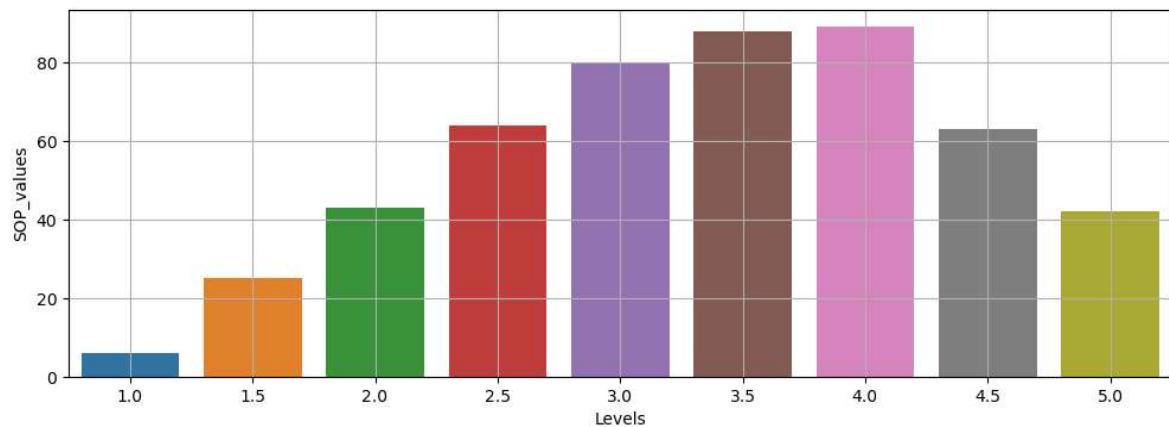
```
In [18]: 1 sop = pd.DataFrame(df['SOP'].value_counts()).reset_index()
```

```
In [19]: 1 sop.rename(columns={"index":"Levels", "SOP":"SOP_values"}, inplace=True)
```

```
In [20]: 1 sop.sort_index
```

```
Out[20]: <bound method DataFrame.sort_index of      Levels  SOP_values
0      4.0        89
1      3.5        88
2      3.0        80
3      2.5        64
4      4.5        63
5      2.0        43
6      5.0        42
7      1.5        25
8      1.0        6>
```

```
In [21]: 1 plt.figure(figsize=(12,4))
2 sns.barplot(data=sop, x='Levels', y='SOP_values')
3 plt.grid()
4 plt.show()
```



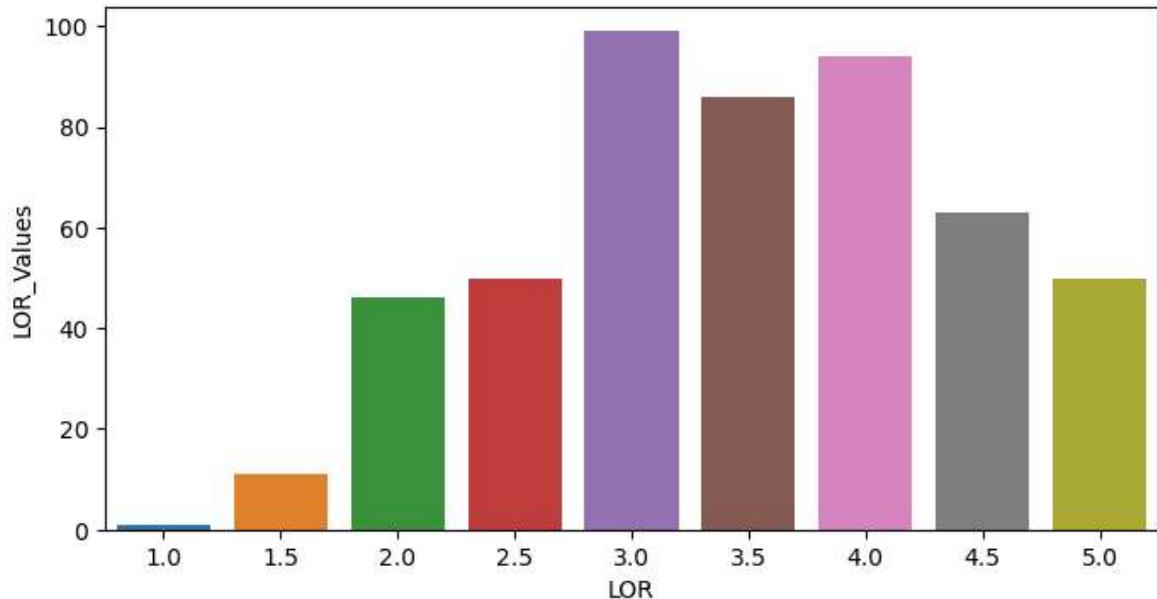
LOR - Letter of Recommendation Strength

```
In [22]: 1 LOR = pd.DataFrame(df['LOR'].value_counts()).reset_index()
2 LOR
```

```
Out[22]:
```

index	LOR
0	3.0 99
1	4.0 94
2	3.5 86
3	4.5 63
4	2.5 50
5	5.0 50
6	2.0 46
7	1.5 11
8	1.0 1

```
In [23]: 1 plt.figure(figsize=(8,4))
2 sns.barplot(data=LOR, x='index', y='LOR ')
3 plt.xlabel("LOR")
4 plt.ylabel("LOR_Values")
5 plt.show()
```



Feature Engineering

Top Rated Chances of Admit Students

```
In [24]: 1 df[df['Chance of Admit '] == df['Chance of Admit '].max()]
```

Out[24]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
24	336	119	5	4.0	3.5	9.80	1	0.97
143	340	120	4	4.5	4.0	9.92	1	0.97
202	340	120	5	4.5	4.5	9.91	1	0.97
203	334	120	5	4.0	5.0	9.87	1	0.97

Lowest Rated chance of Admit of students

```
In [25]: 1 df[df['Chance of Admit '] == df['Chance of Admit '].min()]
```

Out[25]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
92	298	98	2	4.0	3.0	8.03	0	0.34
376	297	96	2	2.5	2.0	7.43	0	0.34

Top CGPA Students

```
In [26]: 1 df[df['CGPA'] == df['CGPA'].max()]
```

Out[26]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
143	340	120	4	4.5	4.0	9.92	1	0.97

Lowest CGPA Students

```
In [27]: 1 df[df['CGPA'] == df['CGPA'].min()]
```

Out[27]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
58	300	99	1	3.0	2.0	6.8	1	0.36

Top GRE scored Students

In [28]: 1 df[df['GRE Score'] == df['GRE Score'].max()]

Out[28]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
25	340	120	5	4.5	4.5	9.60	1	0.94
33	340	114	5	4.0	4.0	9.60	1	0.90
81	340	120	4	5.0	5.0	9.50	1	0.96
84	340	115	5	4.5	4.5	9.45	1	0.94
143	340	120	4	4.5	4.0	9.92	1	0.97
202	340	120	5	4.5	4.5	9.91	1	0.97
284	340	112	4	5.0	4.5	9.66	1	0.94
384	340	113	4	5.0	5.0	9.74	1	0.96
429	340	115	5	5.0	4.5	9.06	1	0.95

Lowest GRE scored Students

In [29]: 1 df[df['GRE Score'] == df['GRE Score'].min()]

Out[29]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
117	290	104	4	2.0	2.5	7.46	0	0.45
377	290	100	1	1.5	2.0	7.56	0	0.47

Top TOEFL Scored Indivisulas

In [30]: 1 df[df['TOEFL Score'] == df['TOEFL Score'].max()]

Out[30]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
25	340	120	5	4.5	4.5	9.60	1	0.94
81	340	120	4	5.0	5.0	9.50	1	0.96
97	331	120	3	4.0	4.0	8.96	1	0.86
143	340	120	4	4.5	4.0	9.92	1	0.97
202	340	120	5	4.5	4.5	9.91	1	0.97
203	334	120	5	4.0	5.0	9.87	1	0.97
212	338	120	4	5.0	5.0	9.66	1	0.95
297	320	120	3	4.0	4.5	9.11	0	0.86
497	330	120	5	4.5	5.0	9.56	1	0.93

Lowest TOEFL Score Students

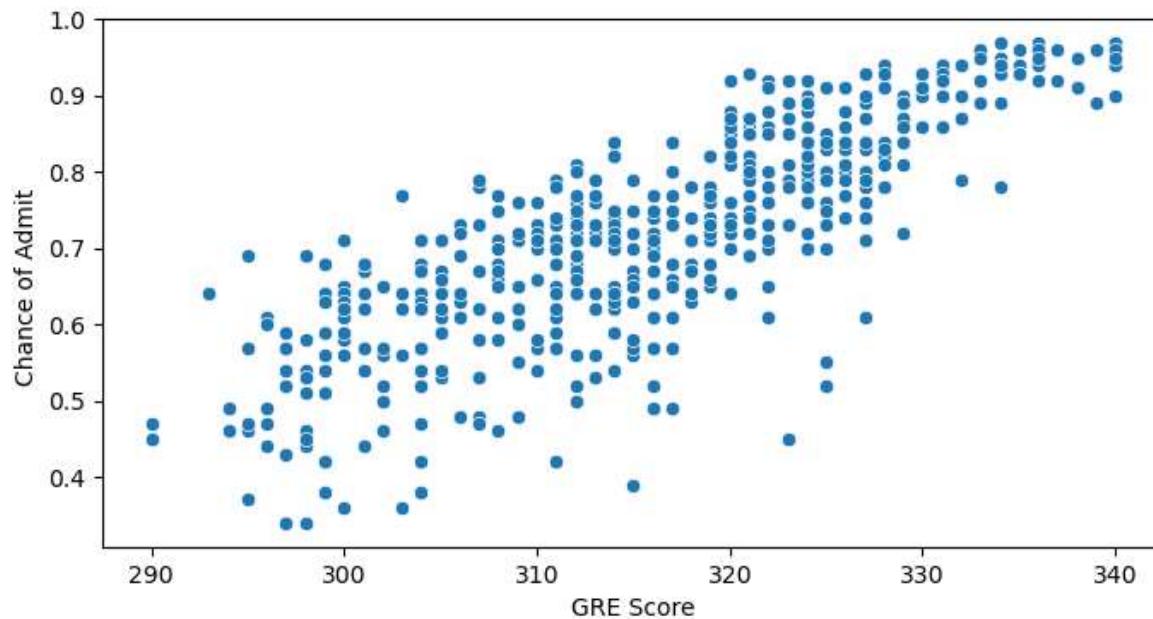
In [31]: 1 df[df['TOEFL Score'] == df['TOEFL Score'].min()]

Out[31]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
368	298	92	1	2.0	2.0	7.88	0	0.51

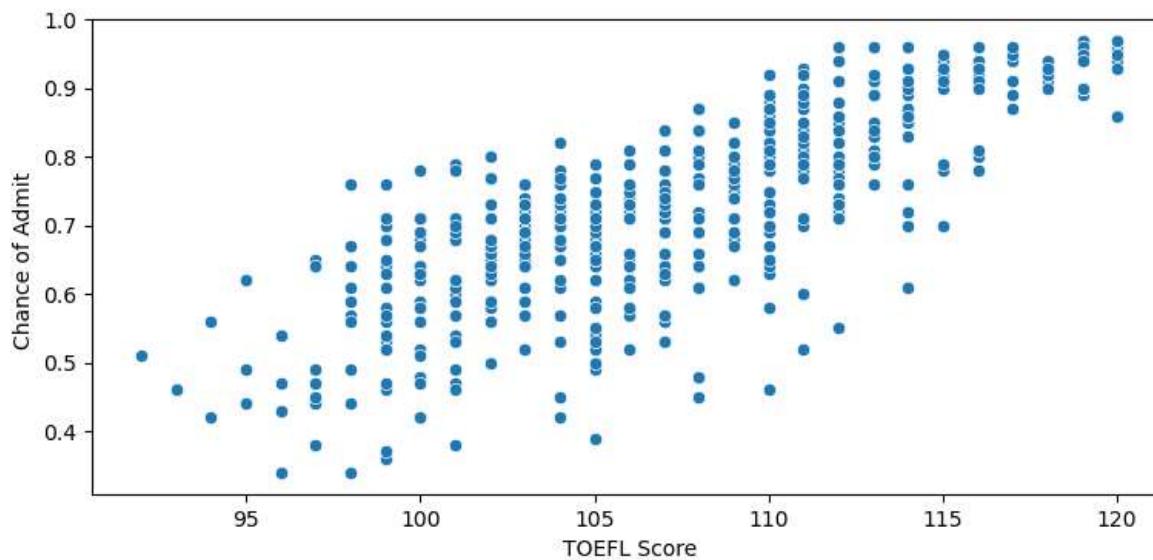
Checking for the Two columns Analysis, Bivariate

```
In [32]:  
1 plt.figure(figsize=(8,4))  
2 sns.scatterplot(data=df, x='GRE Score', y='Chance of Admit ')  
3 plt.show()
```



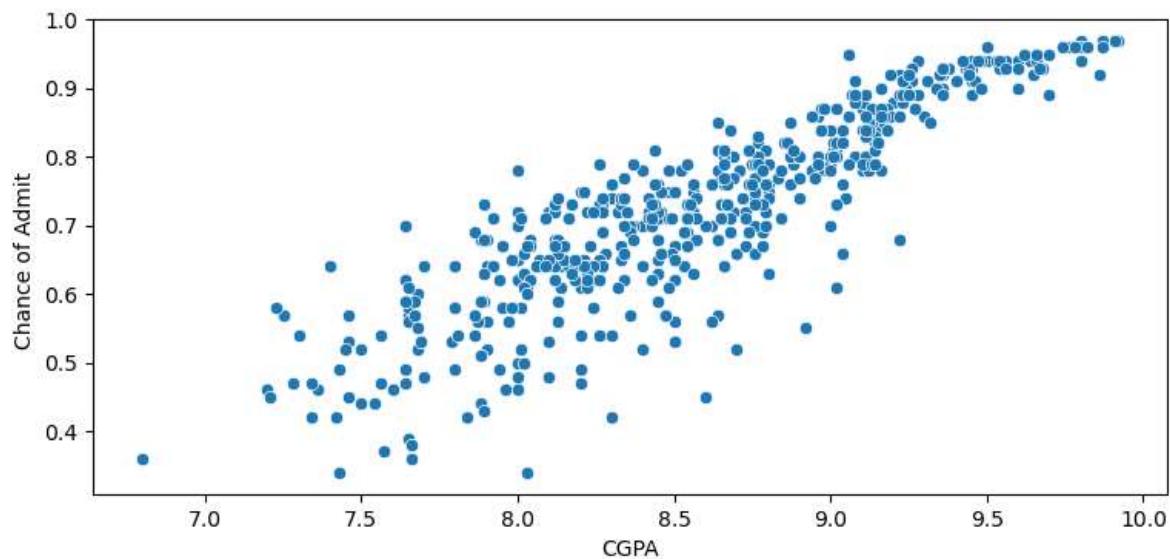
- 1 The scatter plot reveals a nearly linear trend between the variables, indicating a positive correlation with an upward trajectory.

```
In [33]:  
1 plt.figure(figsize=(9,4))  
2 sns.scatterplot(data=df, x='TOEFL Score', y='Chance of Admit ')  
3 plt.show()
```



In [34]:

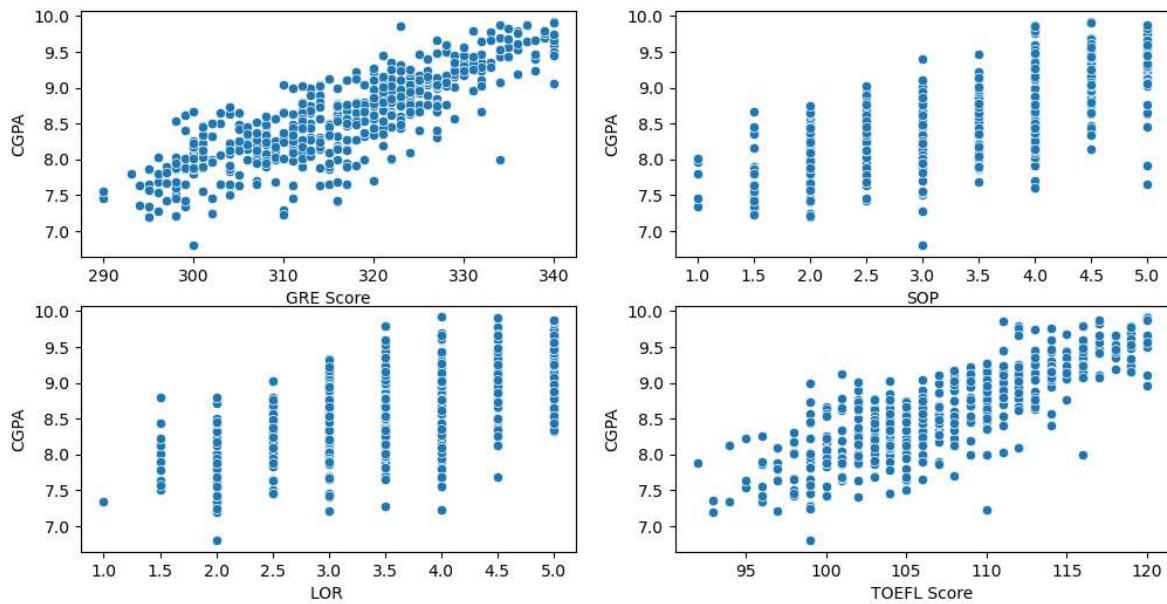
```
1 plt.figure(figsize=(9,4))
2 sns.scatterplot(data=df, x='CGPA', y='Chance of Admit ')
3 plt.show()
```



- 1 The scatter plot above illustrates a discernible linear relationship between CGPA and the likelihood of admission. The upward trajectory indicates that as CGPA increases, there is a corresponding rise in the probability of admission, ranging from significantly low to highly precise values along the linear trend.

In [35]:

```
1 plt.figure(figsize=(12,6))
2 plt.subplot(2,2,1)
3 sns.scatterplot(data=df, x='GRE Score', y='CGPA')
4 plt.subplot(2,2,2)
5 sns.scatterplot(data=df, x='SOP', y='CGPA')
6 plt.subplot(2,2,3)
7 sns.scatterplot(data=df, x='LOR ', y='CGPA')
8 plt.subplot(2,2,4)
9 sns.scatterplot(data=df, x='TOEFL Score', y='CGPA')
10 plt.show()
```



In [36]:

```
1 pd.DataFrame(df.corr())
```

Out[36]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
GRE Score	1.000000	0.827200	0.635376	0.613498	0.524679	0.825878	0.563398	0.810351
TOEFL Score	0.827200	1.000000	0.649799	0.644410	0.541563	0.810574	0.467012	0.792228
University Rating	0.635376	0.649799	1.000000	0.728024	0.608651	0.705254	0.427047	0.690132
SOP	0.613498	0.644410	0.728024	1.000000	0.663707	0.712154	0.408116	0.684137
LOR	0.524679	0.541563	0.608651	0.663707	1.000000	0.637469	0.372526	0.645365
CGPA	0.825878	0.810574	0.705254	0.712154	0.637469	1.000000	0.501311	0.882413
Research	0.563398	0.467012	0.427047	0.408116	0.372526	0.501311	1.000000	0.545871
Chance of Admit	0.810351	0.792228	0.690132	0.684137	0.645365	0.882413	0.545871	1.000000

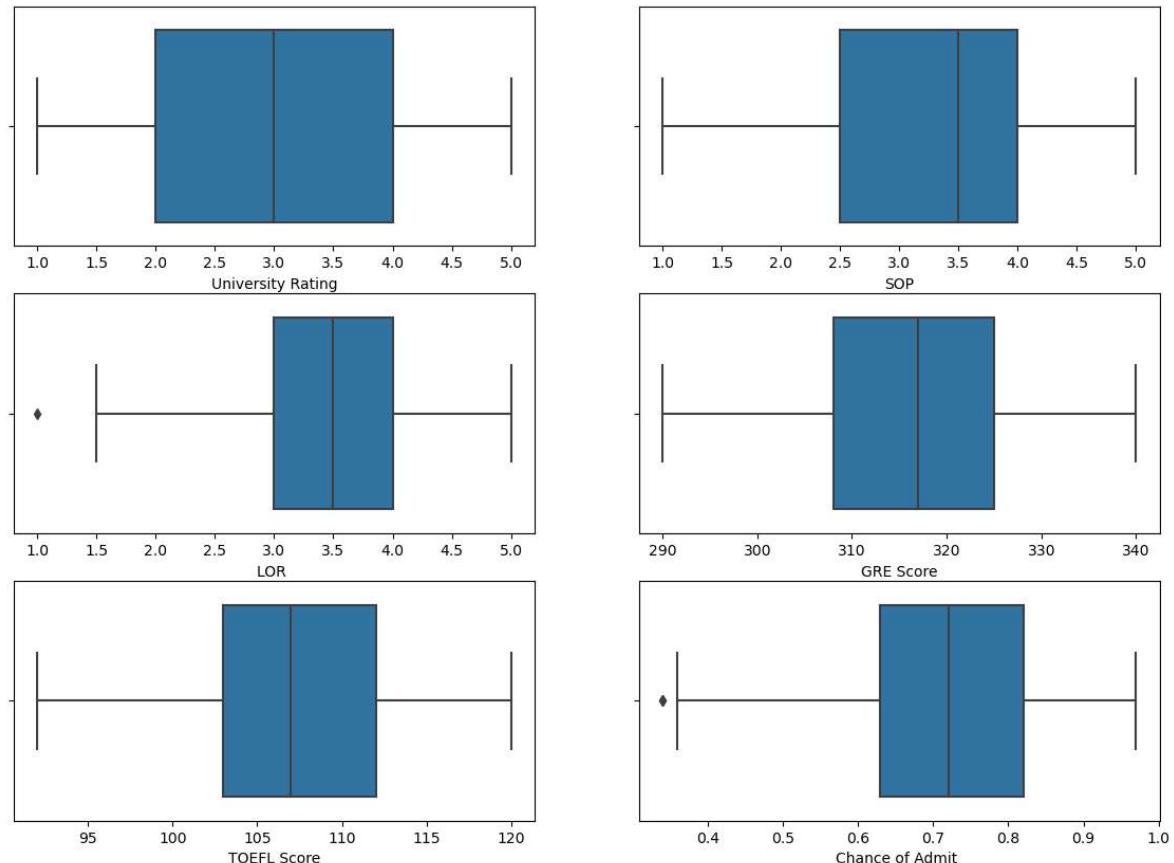
In [37]:

```
1 plt.figure(figsize=(12,5))
2 sns.heatmap(df.corr(), annot=True)
3 plt.show()
```



In [38]:

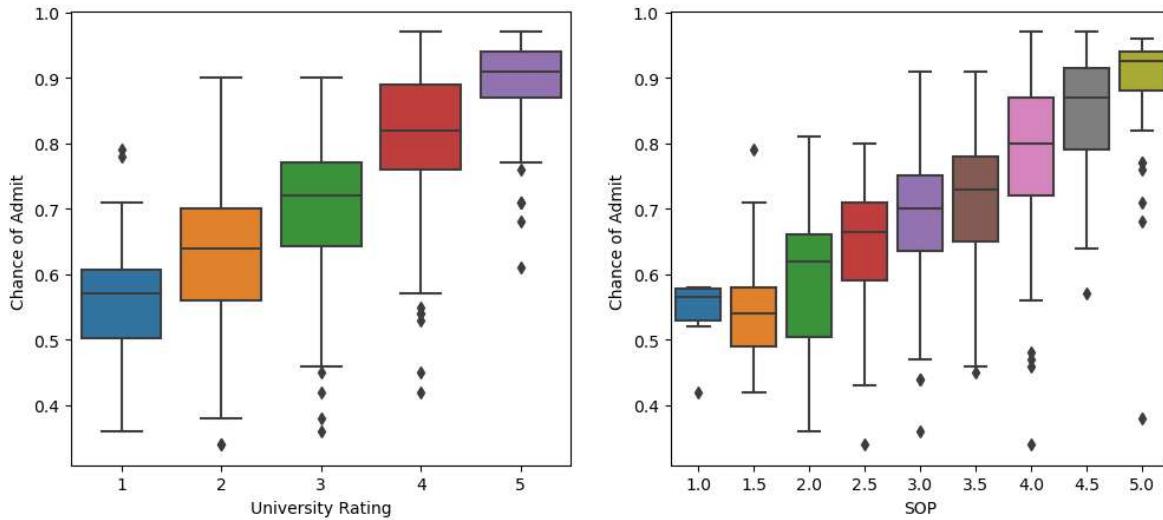
```
1 plt.figure(figsize=(14,10))
2 plt.subplot(3,2,1)
3 sns.boxplot(data=df, x='University Rating')
4 plt.subplot(3,2,2)
5 sns.boxplot(data=df, x='SOP')
6 plt.subplot(3,2,3)
7 sns.boxplot(data=df, x='LOR ')
8 plt.subplot(3,2,4)
9 sns.boxplot(data=df, x='GRE Score')
10 plt.subplot(3,2,5)
11 sns.boxplot(data=df, x='TOEFL Score')
12 plt.subplot(3,2,6)
13 sns.boxplot(data=df, x='Chance of Admit ')
14 plt.show()
```



- Upon examination of each feature, it is evident that the majority exhibit an absence of errors or outliers. This suggests that the data is well-organized, and there are no instances of students significantly deviating from the norm.

In [39]:

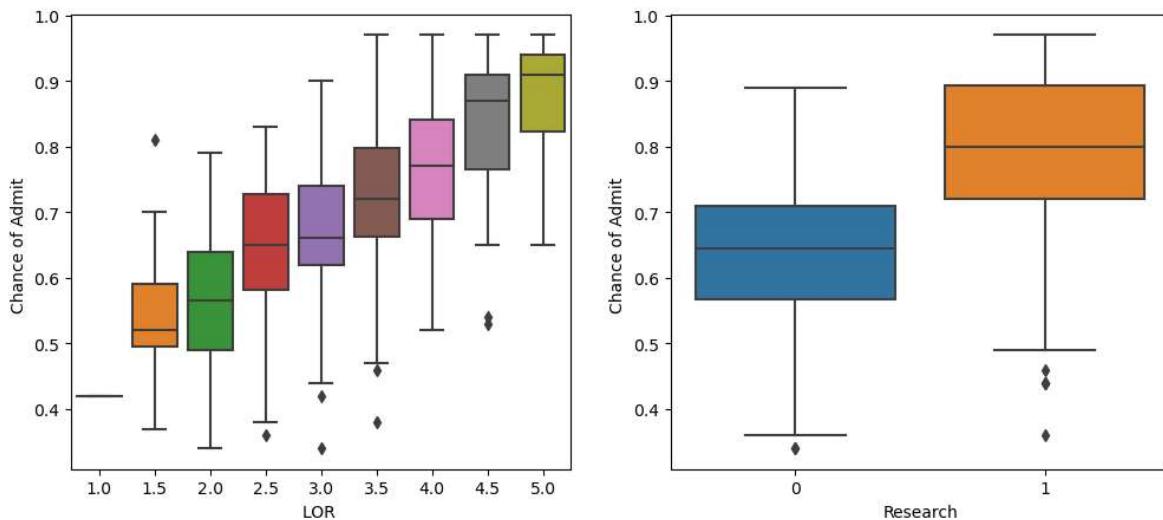
```
1 plt.figure(figsize=(12,5))
2 plt.subplot(1,2,1)
3 sns.boxplot(data=df,x='University Rating', y='Chance of Admit ')
4 plt.subplot(1,2,2)
5 sns.boxplot(data=df,x='SOP', y='Chance of Admit ')
6 plt.show()
```



- During the interpretation of the relationship between the chance of admission and university ranking, the box plot reveals outliers in the ratings of 3, 4, and 5. These outliers, although exceeding standard parameters, may be considered acceptable. Conversely, in the Statement of Purpose (SOP), outliers are also observed, correlating with specific ratings.

In [40]:

```
1 plt.figure(figsize=(12,5))
2 plt.subplot(1,2,1)
3 sns.boxplot(data=df,x='LOR ', y='Chance of Admit ')
4 plt.subplot(1,2,2)
5 sns.boxplot(data=df,x='Research', y='Chance of Admit ')
6 plt.show()
```



1 During the interpretation of the relationship between the chance of admission and LOR, the box plot reveals outliers in the ratings of 3, 3.5, and 4.5. These outliers, although exceeding standard parameters, may be considered acceptable. Conversely, in the Research, outliers are also observed, correlating with specific ratings.

In [41]:

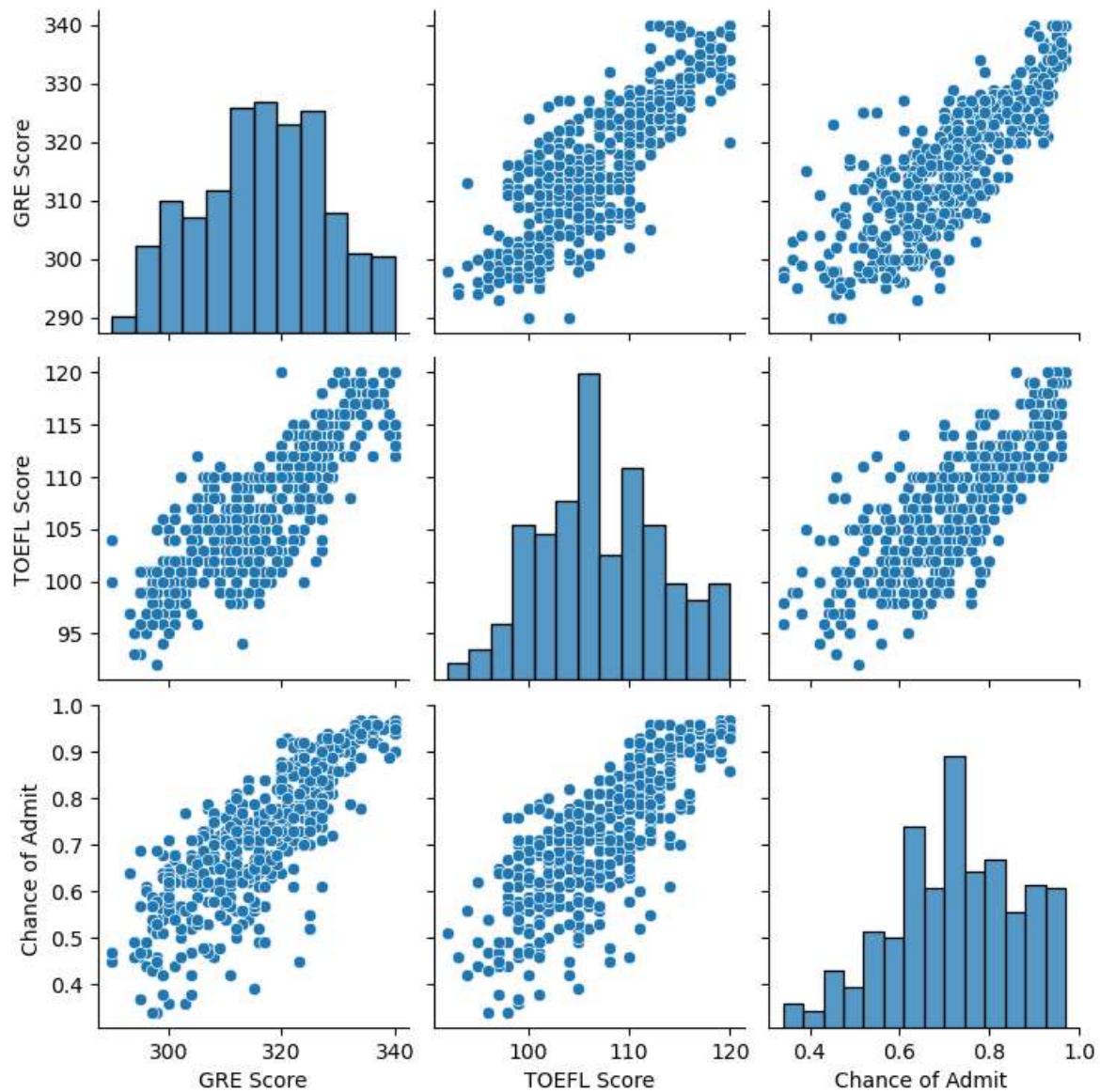
```
1 columns = df[['GRE Score', 'TOEFL Score','Chance of Admit ']]  
2 columns
```

Out[41]:

	GRE Score	TOEFL Score	Chance of Admit
0	337	118	0.92
1	324	107	0.76
2	316	104	0.72
3	322	110	0.80
4	314	103	0.65
...
495	332	108	0.87
496	337	117	0.96
497	330	120	0.93
498	312	103	0.73
499	327	113	0.84

500 rows × 3 columns

```
In [42]: 1 sns.pairplot(columns)
          2 plt.show()
```



```
In [43]: 1 df.duplicated(keep='first').value_counts()
```

Out[43]: False 500
dtype: int64

```
1 # Here we can see there is no duplicated values present in the data set
```

Checking for the Outliers

GRE Score outliers

In [44]:

```
1 q1 = np.quantile(df['GRE Score'], 0.25)
2 q3 = np.quantile(df['GRE Score'], 0.75)
3 print("The 25th percentile", q1, "--", "The 75th percentile",q3)
4 IQR = q3-q1
5 upper_whisker = q3 +(1.5*IQR)
6 Lower_whisker = q1 -(1.5*IQR)
7 print(df[(df['GRE Score'] < Lower_whisker) | (df['GRE Score'] > upper_whisker)])
8 print("No outliers are Present in GRE Score")
```

The 25th percentile 308.0 -- The 75th percentile 325.0
Empty DataFrame
Columns: [GRE Score, TOEFL Score, University Rating, SOP, LOR , CGPA, Research, Chance of Admit]
Index: []
No outliers are Present in GRE Score

TOEFL Score Outliers

In [45]:

```
1 q1 = np.quantile(df['TOEFL Score'], 0.25)
2 q3 = np.quantile(df['TOEFL Score'], 0.75)
3 print("The 25th percentile", q1, "--", "The 75th percentile",q3)
4 IQR = q3-q1
5 upper_whisker = q3 +(1.5*IQR)
6 Lower_whisker = q1 -(1.5*IQR)
7 print(df[(df['TOEFL Score'] < Lower_whisker) | (df['TOEFL Score'] > upper_whisker)])
8 print("No outliers are Present in TOEFL Score")
```

The 25th percentile 103.0 -- The 75th percentile 112.0
Empty DataFrame
Columns: [GRE Score, TOEFL Score, University Rating, SOP, LOR , CGPA, Research, Chance of Admit]
Index: []
No outliers are Present in TOEFL Score

University Rating Outliers

In [46]:

```
1 q1 = np.quantile(df['University Rating'], 0.25)
2 q3 = np.quantile(df['University Rating'], 0.75)
3 print("The 25th percentile", q1, "--", "The 75th percentile", q3)
4 IQR = q3-q1
5 upper_whisker = q3 +(1.5*IQR)
6 Lower_whisker = q1 -(1.5*IQR)
7 df[(df['University Rating'] < Lower_whisker) | (df['University Rating'] > upper_whisker)]
8 print("No outliers are Present in University Rating Score")
```

The 25th percentile 2.0 -- The 75th percentile 4.0
Empty DataFrame
Columns: [GRE Score, TOEFL Score, University Rating, SOP, LOR , CGPA, Research, Chance of Admit]
Index: []
No outliers are Present in University Rating Score

SOP Outliers

In [47]:

```
1 q1 = np.quantile(df['SOP'], 0.25)
2 q3 = np.quantile(df['SOP'], 0.75)
3 print("The 25th percentile", q1, "--", "The 75th percentile", q3)
4 IQR = q3-q1
5 upper_whisker = q3 +(1.5*IQR)
6 Lower_whisker = q1 -(1.5*IQR)
7 df[(df['SOP'] < Lower_whisker) | (df['SOP'] > upper_whisker)]
```

The 25th percentile 2.5 -- The 75th percentile 4.0

Out[47]:

GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
-----------	-------------	-------------------	-----	-----	------	----------	-----------------

LOR Outliers

In [48]:

```
1 q1 = np.quantile(df['LOR '], 0.25)
2 q3 = np.quantile(df['LOR '], 0.75)
3 print("The 25th percentile", q1, "--", "The 75th percentile", q3)
4 IQR = q3-q1
5 upper_whisker = q3 +(1.5*IQR)
6 Lower_whisker = q1 -(1.5*IQR)
7 df[(df['LOR '] < Lower_whisker) | (df['LOR '] > upper_whisker)]
```

The 25th percentile 3.0 -- The 75th percentile 4.0

Out[48]:

GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit	
347	299	94	1	1.0	1.0	7.34	0	0.42

CGPA Outliers

In [49]:

```
1 q1 = np.quantile(df['CGPA'], 0.25)
2 q3 = np.quantile(df['CGPA'], 0.75)
3 print("The 25th percentile", q1, "--", "The 75th percentile", q3)
4 IQR = q3-q1
5 upper_whisker = q3 +(1.5*IQR)
6 Lower_whisker = q1 -(1.5*IQR)
7 df[(df['CGPA'] < Lower_whisker) | (df['CGPA'] > upper_whisker)]
```

The 25th percentile 8.127500000000001 -- The 75th percentile 9.04

Out[49]:

GRE Score TOEFL Score University Rating SOP LOR CGPA Research Chance of Admit

Research Outliers

In [50]:

```
1 q1 = np.quantile(df['Research'], 0.25)
2 q3 = np.quantile(df['Research'], 0.75)
3 print("The 25th percentile", q1, "--", "The 75th percentile", q3)
4 IQR = q3-q1
5 upper_whisker = q3 +(1.5*IQR)
6 Lower_whisker = q1 -(1.5*IQR)
7 df[(df['Research'] < Lower_whisker) | (df['Research'] > upper_whisker)]
```

The 25th percentile 0.0 -- The 75th percentile 1.0

Out[50]:

GRE Score TOEFL Score University Rating SOP LOR CGPA Research Chance of Admit

Chance of Admit

In [51]:

```
1 q1 = np.quantile(df['Chance of Admit'], 0.25)
2 q3 = np.quantile(df['Chance of Admit'], 0.75)
3 print("The 25th percentile", q1, "--", "The 75th percentile", q3)
4 IQR = q3-q1
5 upper_whisker = q3 +(1.5*IQR)
6 Lower_whisker = q1 -(1.5*IQR)
7 df[(df['Chance of Admit'] < Lower_whisker) | (df['Chance of Admit'] >
```

The 25th percentile 0.63 -- The 75th percentile 0.82

Out[51]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
92	298	98	2	4.0	3.0	8.03	0	0.34
376	297	96	2	2.5	2.0	7.43	0	0.34

1 Within the "Chance of Admit" feature, outliers are evident in the dataset, indicating that the likelihood of admission for certain individuals is notably lower compared to their peers. Specifically, students with IDs 92 and 376 may need to exert additional effort to secure admission to prestigious international colleges.

In [52]: 1 df.head()

Out[52]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit	
0	337	118		4	4.5	4.5	9.65	1	0.92
1	324	107		4	4.0	4.5	8.87	1	0.76
2	316	104		3	3.0	3.5	8.00	1	0.72
3	322	110		3	3.5	2.5	8.67	1	0.80
4	314	103		2	2.0	3.0	8.21	0	0.65

In [53]: 1 df.columns

Out[53]: Index(['GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR ', 'CGP A',
'Research', 'Chance of Admit '],
dtype='object')

Data Preparation for Modelling

1 Let us calculate the target variable for prediction in the Linear Regression algorithm. The target variable, representing the chance of admission, will be derived from the dataset. This variable serves as the dependent variable to be predicted based on the selected features.

In [54]: 1 df.head()

Out[54]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit	
0	337	118		4	4.5	4.5	9.65	1	0.92
1	324	107		4	4.0	4.5	8.87	1	0.76
2	316	104		3	3.0	3.5	8.00	1	0.72
3	322	110		3	3.5	2.5	8.67	1	0.80
4	314	103		2	2.0	3.0	8.21	0	0.65

```
In [55]: 1 y = df['Chance of Admit ']
2 y.head()
```

```
Out[55]: 0    0.92
1    0.76
2    0.72
3    0.80
4    0.65
Name: Chance of Admit , dtype: float64
```

1 The provided dataset serves the purpose of training to determine the optimal linear regression line. This process involves establishing a relationship between independent variables and the dependent variable. The independent variables are utilized in the training phase to model and predict the desired outcomes.

```
In [56]: 1 X = df.drop(columns="Chance of Admit ")
2 X.head()
```

```
Out[56]:
```

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research
0	337	118		4	4.5	4.5	9.65
1	324	107		4	4.0	4.5	8.87
2	316	104		3	3.0	3.5	8.00
3	322	110		3	3.5	2.5	8.67
4	314	103		2	2.0	3.0	8.21

```
In [57]: 1 y.shape, X.shape
```

```
Out[57]: ((500,), (500, 7))
```

Trainig and Testing the Independent and Dependent Features

```
In [58]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
```

```
In [59]: 1 X_train.head()
```

```
Out[59]:
```

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research
249	321	111		3	3.5	4.0	8.83
433	316	111		4	4.0	5.0	8.54
19	303	102		3	3.5	3.0	8.50
322	314	107		2	2.5	4.0	8.27
332	308	106		3	3.5	2.5	8.21

```
In [60]: 1 X_test.head()
```

Out[60]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research
361	334	116		4	4.0	3.5	9.54
73	314	108		4	4.5	4.0	9.04
374	315	105		2	2.0	2.5	7.65
155	312	109		3	3.0	3.0	8.69
104	326	112		3	3.5	3.0	9.05

Training shape

```
In [61]: 1 print(X_train.shape)
2 print(y_train.shape)
```

(400, 7)
(400,)

Testing shape of the data

```
In [62]: 1 print(X_test.shape)
2 print(y_test.shape)
```

(100, 7)
(100,)

- 1 Prior to initiating the modeling phase, it is essential to preprocess the data through transformation techniques. This ensures uniformity in the data range and facilitates compatibility with scaling methods such as Min-Max scaling or standard scaling statistics.

```
In [63]: 1 scaling = MinMaxScaler()
2 scaling
```

Out[63]:

└ MinMaxScaler

 MinMaxScaler()

- 1 In this context, the method employed is `fit_transform` to adapt the data for model fitting. This transformation is specifically applied to the `x_test` dataset, aiming to achieve normality in the test data for optimal algorithmic performance.

```
In [64]: 1 X_train = pd.DataFrame(scaling.fit_transform(X_train), columns=X_train.co  
2 X_train
```

Out[64]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	
0	0.62	0.678571		0.50	0.625	0.714286	0.650641	1.0
1	0.52	0.678571		0.75	0.750	1.000000	0.557692	0.0
2	0.26	0.357143		0.50	0.625	0.428571	0.544872	0.0
3	0.48	0.535714		0.25	0.375	0.714286	0.471154	0.0
4	0.36	0.500000		0.50	0.625	0.285714	0.451923	1.0
...
395	0.78	0.678571		0.75	0.875	0.857143	0.762821	1.0
396	0.32	0.464286		0.25	0.375	0.428571	0.455128	1.0
397	0.24	0.250000		0.00	0.250	0.142857	0.144231	0.0
398	0.38	0.464286		0.25	0.375	0.714286	0.282051	0.0
399	0.48	0.500000		0.25	0.750	0.571429	0.464744	0.0

400 rows × 7 columns

```
In [65]: 1 X_test = pd.DataFrame(scaling.transform(X_test), columns=X_test.columns)  
2 X_test
```

Out[65]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	
0	0.88	0.857143		0.75	0.750	0.571429	0.878205	1.0
1	0.48	0.571429		0.75	0.875	0.714286	0.717949	1.0
2	0.50	0.464286		0.25	0.250	0.285714	0.272436	0.0
3	0.44	0.607143		0.50	0.500	0.428571	0.605769	0.0
4	0.72	0.714286		0.50	0.625	0.428571	0.721154	1.0
...
95	0.18	0.071429		0.00	0.000	-0.142857	0.173077	0.0
96	0.50	0.500000		0.50	0.875	0.571429	0.519231	0.0
97	0.78	0.785714		0.25	0.250	0.714286	0.564103	1.0
98	0.56	0.642857		0.00	0.375	0.571429	0.557692	1.0
99	0.48	0.464286		0.50	0.625	0.285714	0.480769	0.0

100 rows × 7 columns

```
In [66]: 1 print("Training shape: " , X_train.shape)
2 print("Testing shape :" , X_test.shape)
```

Training shape: (400, 7)
Testing shape : (100, 7)

Preparing the Linear Regression Model

```
In [67]: 1 Model = LinearRegression()
2 Model
```

```
Out[67]: ▾ LinearRegression
LinearRegression()
```

```
In [68]: 1 Model.fit(X_train, y_train)
```

```
Out[68]: ▾ LinearRegression
LinearRegression()
```

Feature name with the Coefficient

```
In [69]: 1 pd.DataFrame({'Feature': X_train.columns, 'Coefficient': Model.coef_})
```

```
Out[69]:
```

	Feature	Coefficient
0	GRE Score	0.121722
1	TOEFL Score	0.083884
2	University Rating	0.010275
3	SOP	0.007255
4	LOR	0.060333
5	CGPA	0.351085
6	Research	0.024027

The intercept of the best Line from Linear Regression (W0)

```
In [70]: 1 Model.intercept_
```

```
Out[70]: 0.35558405022527056
```

Train and Test Performance Checking

Score of the Model / The Performance

```
In [71]: 1 r2 = Model.score(X_train, y_train)  
2 r2
```

```
Out[71]: 0.8210671369321554
```

```
In [72]: 1 r2 = Model.score(X_test, y_test)  
2 r2
```

```
Out[72]: 0.8188432567829629
```

```
In [73]: 1 y_test[:10]
```

```
Out[73]: 361    0.93  
73     0.84  
374    0.39  
155    0.77  
104    0.74  
394    0.89  
377    0.47  
124    0.57  
68     0.68  
450    0.82  
Name: Chance of Admit , dtype: float64
```

```
In [74]: 1 y_pred = Model.predict(X_test)  
2 y_pred[:10]
```

```
Out[74]: array([0.91457473, 0.79518127, 0.57265986, 0.70736968, 0.81588282,  
   0.86206561, 0.47459746, 0.64850923, 0.82378728, 0.80741498])
```

Performance with the R2 score

```
In [75]: 1 r2_score(y_test, y_pred)
```

```
Out[75]: 0.8188432567829629
```

Checking for the MSE and RMSE

MSE

```
1 Lets Analyse the Results for the MSE (Mean Sqaure Error)
```

- | | |
|---|--|
| 2 | Mean Squared Error (MSE) gauges how much error exists in a statistical model. It calculates the average squared difference between the predicted and actual values. If the model is perfect with no errors, the MSE is zero. As model errors increase, the MSE value also rises. |
| 3 | Actual training - Predicted of traning data set and squaring them will give MSE |

```
In [76]: 1 np.square(y_train - Model.predict(X_train)).mean()
```

Out[76]: 0.0035265554784557574

Using with the libraries MSE

```
In [77]: 1 mean_squared_error(y_test, y_pred)
```

Out[77]: 0.003704655398788409

MAE (Mean Absolute Error)

- | | |
|---|---|
| 1 | Absolute error measures how far off a prediction is from the actual value. Mean Absolute Error (MAE) calculates the average of these differences across all predictions and observations, providing an overall measure of the magnitude of errors for the entire group. |
|---|---|

```
In [78]: 1 np.abs( y_train - Model.predict(X_train)).mean()
```

Out[78]: 0.042533340611643135

Using with the libraries MAE (Mean Absolute Error)

```
In [79]: 1 mean_absolute_error(y_test, y_pred)
```

Out[79]: 0.04272265427705364

RMSE (Root Mean Sqaure Error)

- | | |
|---|--|
| 1 | The Root Mean Squared Error (RMSE) is a key measure for evaluating how well a regression model predicts values. It calculates the average difference between the predicted values from the model and the actual values. RMSE offers insight into the accuracy of the model's predictions for the target value. |
| 2 | |
| 3 | formula for the RMSE |
| 4 | formula : sqrt(sum(pred - actual)**2)/n |

5 | where n: is the no.of rows

```
In [80]: 1 n = y_test.shape  
2 n[0]
```

Out[80]: 100

```
In [81]: 1 np.sqrt(((np.sum(y_pred - y_test))**2)/n[0])
```

Out[81]: 0.05453623717661262

Actual y_test value and Predicted y Comparison

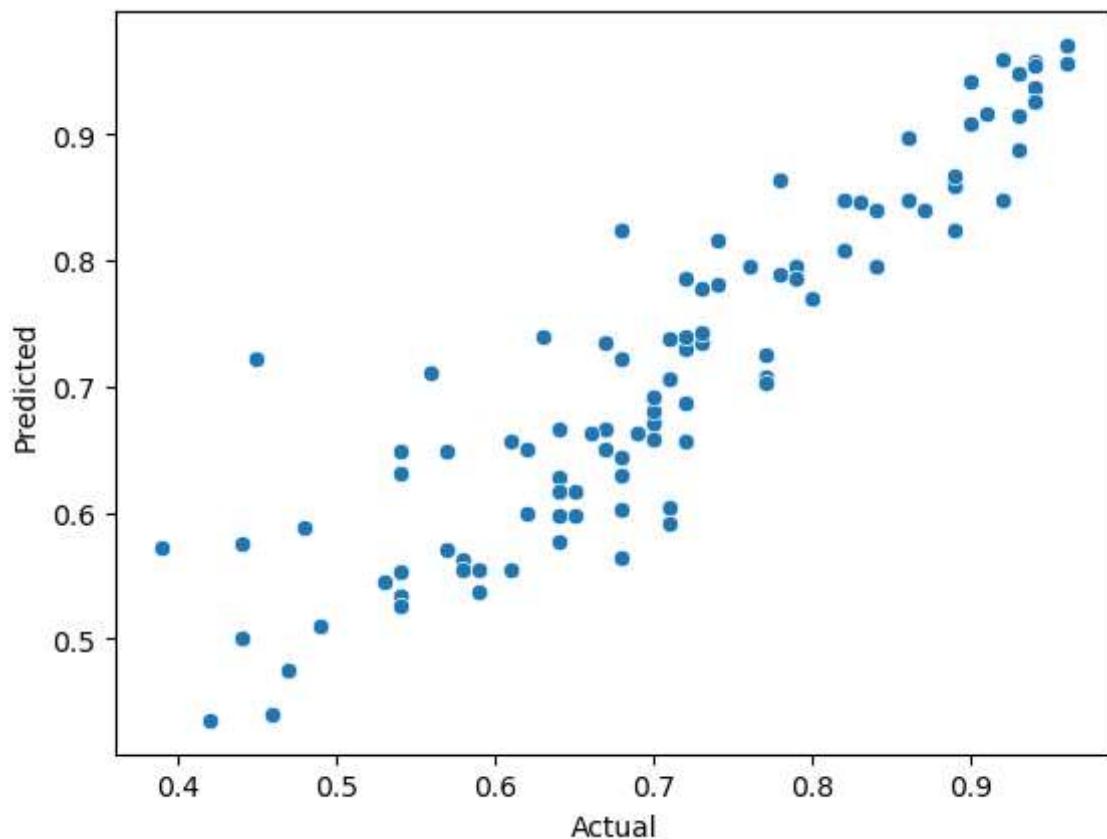
```
In [82]: 1 test1 = pd.DataFrame({  
2     "Actual":y_test,  
3     "Predicted":y_pred  
4 })  
5 test1[:10]
```

Out[82]:

	Actual	Predicted
361	0.93	0.914575
73	0.84	0.795181
374	0.39	0.572660
155	0.77	0.707370
104	0.74	0.815883
394	0.89	0.862066
377	0.47	0.474597
124	0.57	0.648509
68	0.68	0.823787
450	0.82	0.807415

1 | In this phase, the model was trained using the dataset, aligning itself with the dependent values. Utilizing the Linear Regression algorithm, the model achieved an accuracy of approximately 81%. Subsequently, actual predicted values were obtained for specific datasets corresponding to the dependent variables/features.

```
In [83]: 1 sns.scatterplot(data=test1, x="Actual", y="Predicted")
          2 plt.show()
```



- 1 The scatter plot depicting the trend between actual and predicted values for both numerical variables reveals a clear linear pattern. This observation indicates a conformity to linear regression, as the data points align themselves along a linear trend.

Normality of Residuals

Checking for the Error

```
In [84]: 1 Error = test1["Actual"] - test1["Predicted"]
          2 Error[:10]
```

```
Out[84]: 361    0.015425
         73     0.044819
        374   -0.182660
       155    0.062630
      104   -0.075883
      394    0.027934
      377   -0.004597
      124   -0.078509
       68   -0.143787
      450    0.012585
      dtype: float64
```

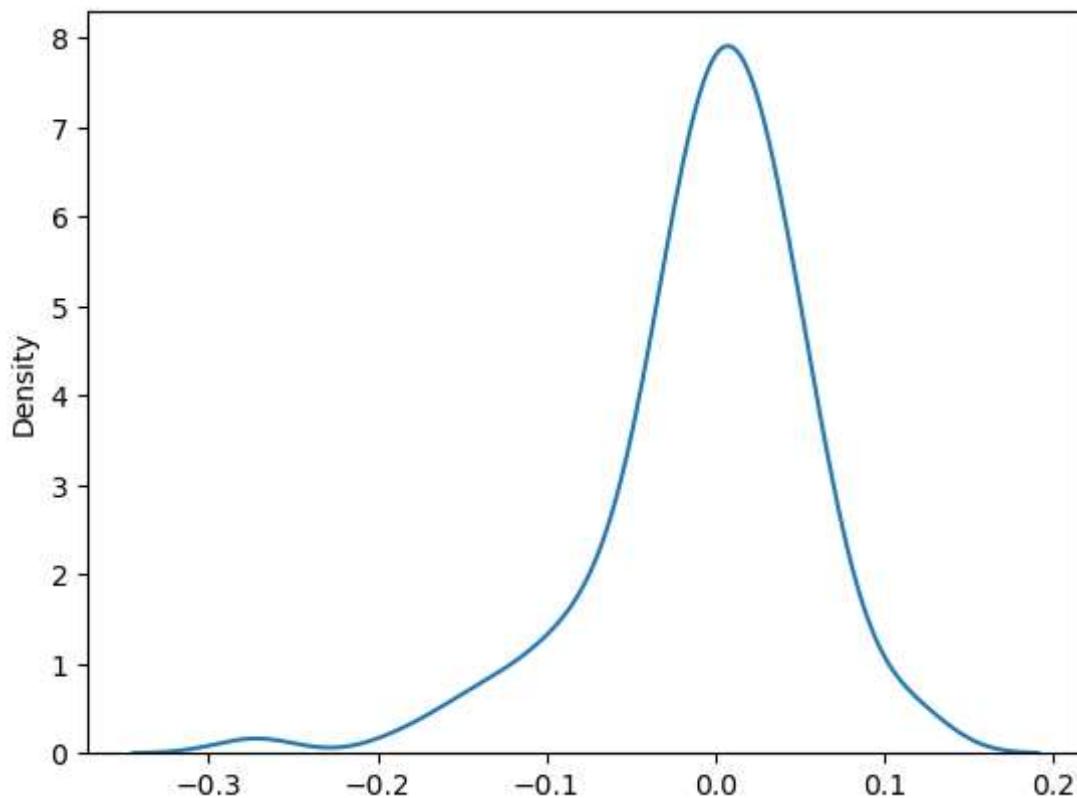
The mean of residuals is nearly zero

```
In [85]: 1 Error.mean()
```

```
Out[85]: -0.005453623717661262
```

- Upon examining the residuals, representing the differences between actual and predicted values, it is observed that their mean is approximately zero. This indicates favorable performance of the residuals, suggesting a cancellation effect and contributing to the normality of the data.

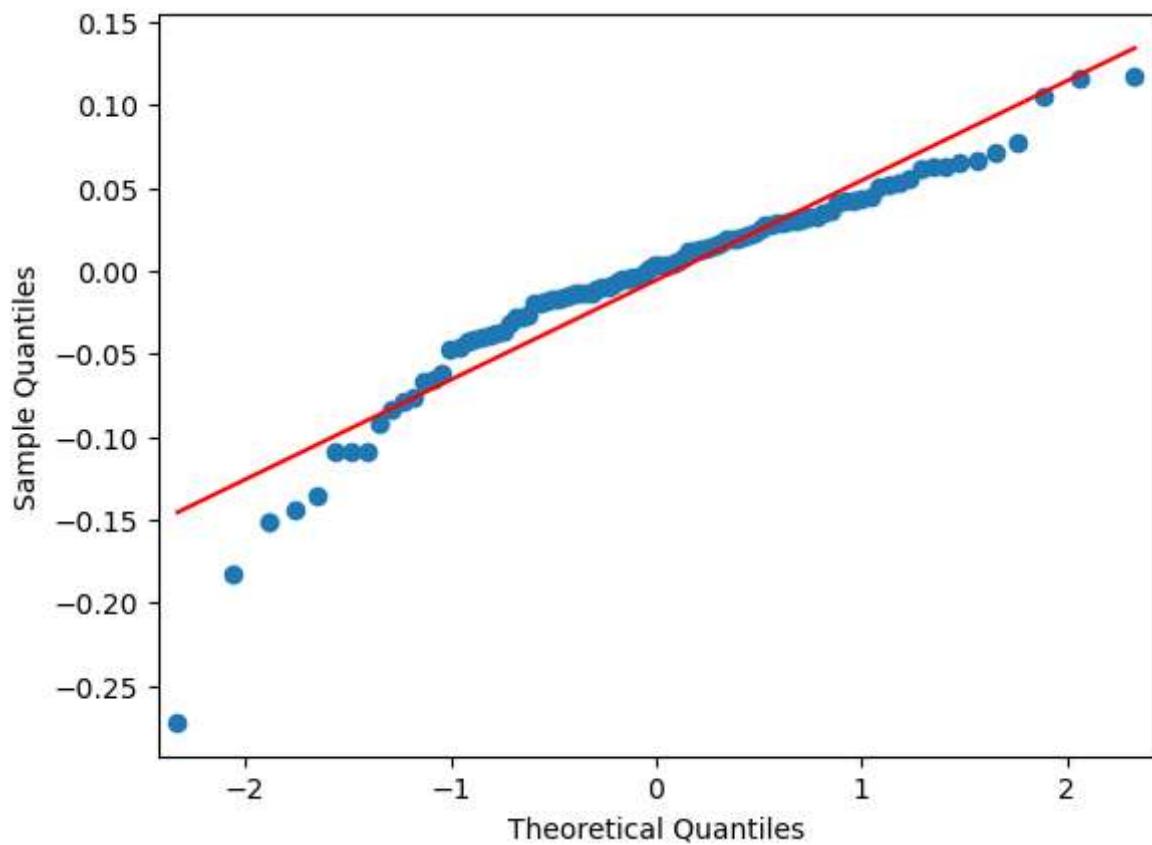
```
In [86]: 1 sns.kdeplot(Error)
          2 plt.show()
```



- Upon scrutinizing the residuals through a kernel density plot, it is evident that the distribution is somewhat skewed, with a mean close to zero. This implies that individuals with relatively lower scores contribute to the observed skewness, despite the overall mean being centered around zero.

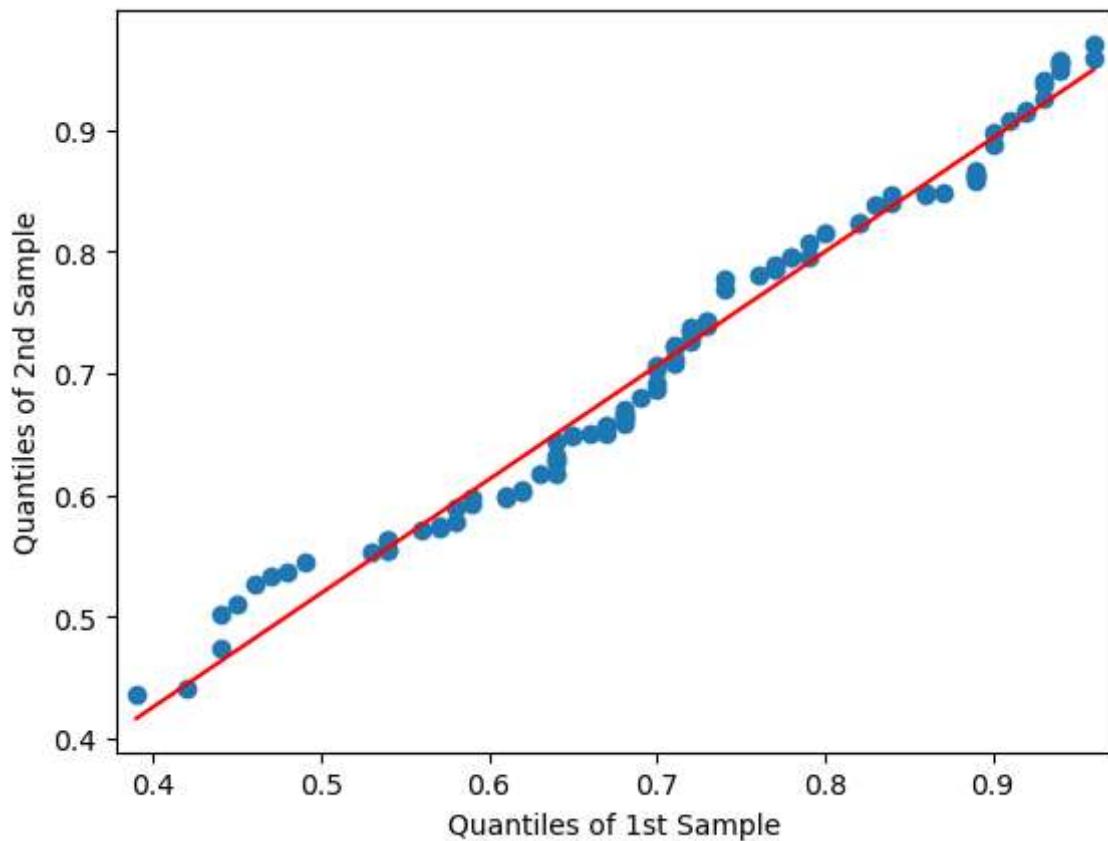
Checking with the QQ plot, for the Regression of fit of line

```
In [87]: 1 sm.qqplot(Error, line="r")
2 plt.show()
```



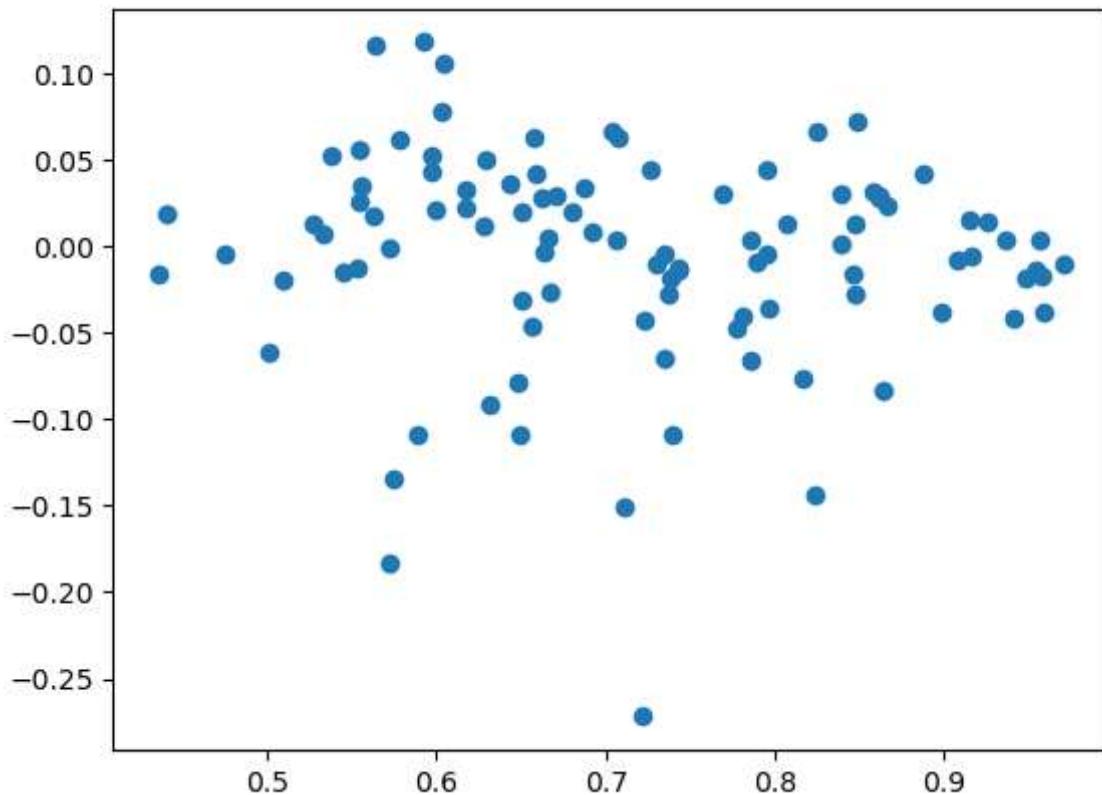
QQ Plot, with Actual and Predicted Line with the best fit of regression

```
In [88]: 1 sm.qqplot_2samples(y_test, y_pred, line='r')
2 plt.show()
```



Checking for the Homoscedasticity

```
In [89]: 1 plt.scatter(y_pred, Error)
          2 plt.show()
```



- 1 The mentioned plot illustrates the dispersion between predicted and actual Errors, indicating Homoscedasticity. However, the data points exhibit significant scattering without a discernible trend, suggesting a notable bias across all values.

Let's Check the Adjusted R2 Score

```
1 Formula : 1 - (1- r2)*(n-1)/(n-d-1)
2 where:
3 r2 is the Score
4 n is the no of Rows
5 d is the no.features
```

```
In [90]: 1 n , d =X_train.shape
2 print("Number of Rows:",n,"Features",d)
3 Numerator = (1-r2)*(n-1)
4 Denominator = (n-d-1)
5 1 - (Numerator/Denominator)
```

Number of Rows: 400 Features 7

Out[90]: 0.8156083149398015

Multicollinearity checking by VIF score

```
In [91]: 1 X_multicheck =pd.DataFrame(X_train, columns=X_train.columns)
2 vif = pd.DataFrame()
3 vif[ 'Features' ] = X_multicheck.columns
4 vif[ 'VIF' ] = [variance_inflation_factor(X_multicheck.values, i) for i in
5 vif[ 'VIF' ] = round(vif[ 'VIF' ], 2)
6 vif = vif.sort_values(by = "VIF", ascending = False)
7 vif
```

Out[91]:

	Features	VIF
5	CGPA	39.76
0	GRE Score	31.20
1	TOEFL Score	26.76
3	SOP	18.57
4	LOR	11.01
2	University Rating	10.95
6	Research	3.36

- 1 The CGPA variable demonstrates a considerably higher VIF score compared to others, indicating a significant level of multicollinearity. Considering the problem statement's criterion of VIF not exceeding 5, the observed multicollinearity necessitates the removal of such variables. This elimination aims to enhance predictive accuracy by mitigating redundant information.

```
In [92]: 1 X_multicheck.drop(columns={"CGPA"}, inplace=True)
```

```
In [93]: 1 X_multicheck.head()
```

Out[93]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	Research
0	0.62	0.678571		0.50	0.625	0.714286
1	0.52	0.678571		0.75	0.750	1.000000
2	0.26	0.357143		0.50	0.625	0.428571
3	0.48	0.535714		0.25	0.375	0.714286
4	0.36	0.500000		0.50	0.625	0.285714

```
In [94]: 1 vif = pd.DataFrame()
```

```
2 vif['Features'] = X_multicheck.columns
3 vif['VIF'] = [variance_inflation_factor(X_multicheck.values, i) for i in
4 vif['VIF'] = round(vif['VIF'], 2)
5 vif = vif.sort_values(by = "VIF", ascending = False)
6 vif
```

Out[94]:

	Features	VIF
0	GRE Score	24.83
1	TOEFL Score	24.22
3	SOP	17.26
2	University Rating	10.90
4	LOR	10.15
5	Research	3.36

- 1 The GRE Score variable demonstrates a considerably higher VIF score compared to others, indicating a significant level of multicollinearity. Considering the problem statement's criterion of VIF not exceeding 5, the observed multicollinearity necessitates the removal of such variables. This elimination aims to enhance predictive accuracy by mitigating redundant information.

```
In [95]: 1 X_multicheck.drop(columns={"GRE Score"}, inplace=True)
```

```
In [96]: 1 vif = pd.DataFrame()
2 vif['Features'] = X_multicheck.columns
3 vif['VIF'] = [variance_inflation_factor(X_multicheck.values, i) for i in
4 vif['VIF'] = round(vif['VIF'], 2)
5 vif = vif.sort_values(by = "VIF", ascending = False)
6 vif
```

Out[96]:

	Features	VIF
2	SOP	17.07
0	TOEFL Score	12.73
1	University Rating	10.79
3	LOR	10.09
4	Research	2.99

- 1 The SOP (Statement of Purpose) variable demonstrates a considerably higher VIF score compared to others, indicating a significant level of multicollinearity. Considering the problem statement's criterion of VIF not exceeding 5, the observed multicollinearity necessitates the removal of such variables. This elimination aims to enhance predictive accuracy by mitigating redundant information.

```
In [97]: 1 X_multicheck.drop(columns={"SOP"}, inplace=True)
```

```
In [98]: 1 vif = pd.DataFrame()
2 vif['Features'] = X_multicheck.columns
3 vif['VIF'] = [variance_inflation_factor(X_multicheck.values, i) for i in
4 vif['VIF'] = round(vif['VIF'], 2)
5 vif = vif.sort_values(by = "VIF", ascending = False)
6 vif
```

Out[98]:

	Features	VIF
0	TOEFL Score	10.51
1	University Rating	9.33
2	LOR	8.17
3	Research	2.98

- 1 The TOEFL Score variable demonstrates a considerably higher VIF score compared to others, indicating a significant level of multicollinearity. Considering the problem statement's criterion of VIF not exceeding 5, the observed multicollinearity necessitates the removal of such variables. This elimination aims to enhance predictive accuracy by mitigating redundant information.

```
In [99]: 1 X_multicheck.drop(columns={"TOEFL Score"}, inplace=True)
```

```
In [100]: 1 vif = pd.DataFrame()
2 vif['Features'] = X_multicheck.columns
3 vif['VIF'] = [variance_inflation_factor(X_multicheck.values, i) for i in
4 vif['VIF'] = round(vif['VIF'], 2)
5 vif = vif.sort_values(by = "VIF", ascending = False)
6 vif
```

Out[100]:

	Features	VIF
0	University Rating	7.19
1	LOR	6.49
2	Research	2.77

1 The University Rating variable demonstrates a considerably higher VIF score compared to others, indicating a significant level of multicollinearity. Considering the problem statement's criterion of VIF not exceeding 5, the observed multicollinearity necessitates the removal of such variables. This elimination aims to enhance predictive accuracy by mitigating redundant information.

```
In [101]: 1 X_multicheck.drop(columns={"University Rating"}, inplace=True)
```

```
In [102]: 1 vif = pd.DataFrame()
2 vif['Features'] = X_multicheck.columns
3 vif['VIF'] = [variance_inflation_factor(X_multicheck.values, i) for i in
4 vif['VIF'] = round(vif['VIF'], 2)
5 vif = vif.sort_values(by = "VIF", ascending = False)
6 vif
```

Out[102]:

	Features	VIF
0	LOR	2.44
1	Research	2.44

Checking Adjusted R2 for No Multicollinearity or less than 5 VIF Score

```
In [103]: 1 n,d = X_multicheck.shape
```

```
In [104]: 1 n , d =X_multicheck.shape  
2 print("Number of Rows:",n,"Features",d)  
3 Numerator = (1-r2)*(n-1)  
4 Denominator = (n-d-1)  
5 1 - (Numerator/Denominator)
```

Number of Rows: 400 Features 2

Out[104]: 0.8179306283536579

1 Previously, the performance, as indicated by the Adjusted R-squared, stood at 81.56. Following the removal of features based on VIF (Variance Inflation Factor), several features were dropped, resulting in a slight increase in performance to 81.79 in the Adjusted R-squared.

Let's Analyse the Stats Summary for the Regression Data using OLS

Here we are adding the constant just because in the OLS (Ordinary least square) do not add the constant in the equation

```
In [105]: 1 x_sm = sm.add_constant(X_train)
```

```
In [106]: 1 y_train1 = y_train.values  
2 y_train1 = y_train1.reshape(-1,1)  
3 y_train1.shape
```

Out[106]: (400, 1)

```
In [107]: 1 model = sm.OLS(y_train1,x_sm)
2 result = model.fit()
3 print(result.summary())
```

OLS Regression Results

Dep. Variable:	y	R-squared:	0.8		
21					
Model:	OLS	Adj. R-squared:	0.8		
18					
Method:	Least Squares	F-statistic:	25		
7.0					
Date:	Tue, 13 Feb 2024	Prob (F-statistic):	3.41e-1		
42					
Time:	14:45:36	Log-Likelihood:	561.		
91					
No. Observations:	400	AIC:	-110		
8.					
Df Residuals:	392	BIC:	-107		
6.					
Df Model:	7				
Covariance Type:	nonrobust				
<hr/>					
<hr/>					
	coef	std err	t	P> t	[0.025
0.975]					
<hr/>					
<hr/>					
const	0.3556	0.010	36.366	0.000	0.336
0.375					
GRE Score	0.1217	0.029	4.196	0.000	0.065
0.179					
TOEFL Score	0.0839	0.026	3.174	0.002	0.032
0.136					
University Rating	0.0103	0.017	0.611	0.541	-0.023
0.043					
SOP	0.0073	0.020	0.357	0.721	-0.033
0.047					
LOR	0.0603	0.016	3.761	0.000	0.029
0.092					
CGPA	0.3511	0.034	10.444	0.000	0.285
0.417					
Research	0.0240	0.007	3.231	0.001	0.009
0.039					
<hr/>					
<hr/>					
Omnibus:	86.232	Durbin-Watson:	2.0		
50					
Prob(Omnibus):	0.000	Jarque-Bera (JB):	190.0		
99					
Skew:	-1.107	Prob(JB):	5.25e-		
42					
Kurtosis:	5.551	Cond. No.	2		
3.4					
<hr/>					
<hr/>					

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Let's Analyse for the better Performance of the model with Polynomial Features

In [108]: 1 df[:2]

Out[108]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit	
0	337	118		4	4.5	4.5	9.65	1	0.92
1	324	107		4	4.0	4.5	8.87	1	0.76

In [109]: 1 X_train.shape

Out[109]: (400, 7)

In [110]: 1 X_test.shape

Out[110]: (100, 7)

In [111]: 1 poly = PolynomialFeatures(degree=2, include_bias=False)
2 poly

Out[111]:

```
PolynomialFeatures  
PolynomialFeatures(include_bias=False)
```

In [112]: 1 model = LinearRegression()
2 model

Out[112]:

```
LinearRegression  
LinearRegression()
```

In [113]: 1 Xtrain = poly.fit_transform(X_train)
2 Xtrain.shape

Out[113]: (400, 35)

In [114]: 1 Xtest = poly.transform(X_test)
2 Xtest.shape

Out[114]: (100, 35)

```
In [115]: 1 model.fit(Xtrain, y_train)
```

```
Out[115]: ▾ LinearRegression  
LinearRegression()
```

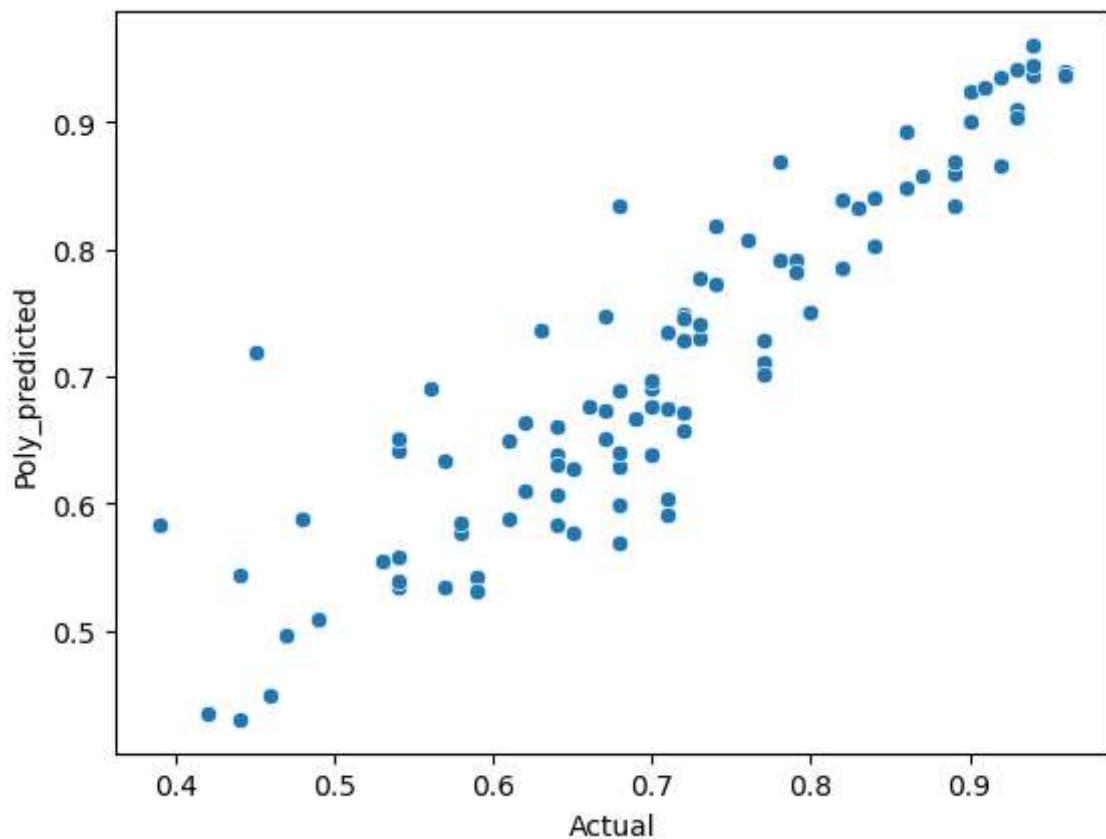
```
In [116]: 1 ypoly_predition = model.predict(Xtest)
```

```
In [117]: 1 polypredicted = pd.DataFrame({  
2     "Actual": y_test,  
3     "Poly_predicted" : model.predict(Xtest)  
4 })  
5 polypredicted[:10]
```

```
Out[117]:
```

	Actual	Poly_predicted
361	0.93	0.910068
73	0.84	0.802416
374	0.39	0.583174
155	0.77	0.710472
104	0.74	0.819335
394	0.89	0.861433
377	0.47	0.497188
124	0.57	0.634287
68	0.68	0.834664
450	0.82	0.785685

```
In [118]: 1 sns.scatterplot(data=polypredicted, x="Actual", y="Poly_predicted")
2 plt.show()
```



Polynomial Prediction Performcae

```
In [119]: 1 model.score(Xtest, y_test)
```

Out[119]: 0.826511555264336

1 Following the implementation of Linear Regression, a score of 0.82106 was attained. Upon incorporating Polynomial Regression, the score marginally increased to 0.8265. Consequently, the difference in performance between Linear Regression and Polynomial Regression is negligible.

With the R2 Score Polynomial Rpediction

```
In [120]: 1 r2_score(y_test, ypoly_prediction)
```

Out[120]: 0.826511555264336

```
1 Following the implementation of Linear Regression, a r2_score of 0.8188  
was attained. Upon incorporating Polynomial Regression, the score  
marginally increased to 0.8265. Consequently, the difference in  
performance between Linear Regression and Polynomial Regression is  
incremented.
```

Adjusted Rscore

```
In [121]: 1 n , d =Xtrain.shape  
2 print("Number of Rows:",n,"Features",d)  
3 Numerator = (1-r2)*(n-1)  
4 Denominator = (n-d-1)  
5 1 - (Numerator/Denominator)
```

Number of Rows: 400 Features 35

Out[121]: 0.8014243391659401

```
1 The inclusion of polynomial features has resulted in a decrease in the  
Adjusted R-squared score. This suggests that the model's introduction  
of additional complexity through polynomial features has weakened its  
overall strength, leading to a reduction in the Adjusted R-squared  
score.
```

Residuals

```
In [122]: 1 polyerror = y_test - ypoly_predition  
2 polyerror[:10]
```

```
Out[122]: 361    0.019932  
73     0.037584  
374    -0.193174  
155    0.059528  
104    -0.079335  
394    0.028567  
377    -0.027188  
124    -0.064287  
68     -0.154664  
450    0.034315  
Name: Chance of Admit , dtype: float64
```

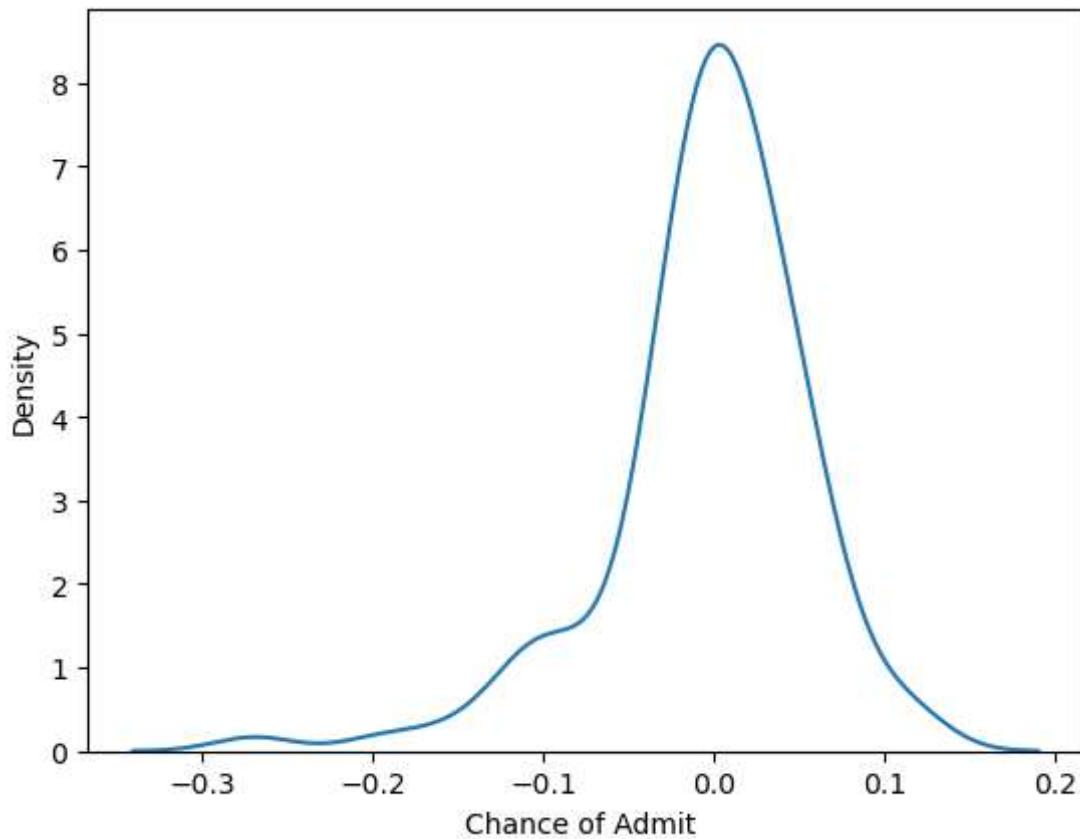
Checking for the Mean of residuals

```
In [123]: 1 np.mean(polyerror)
```

Out[123]: -0.004070822574379355

Normality of Residuals

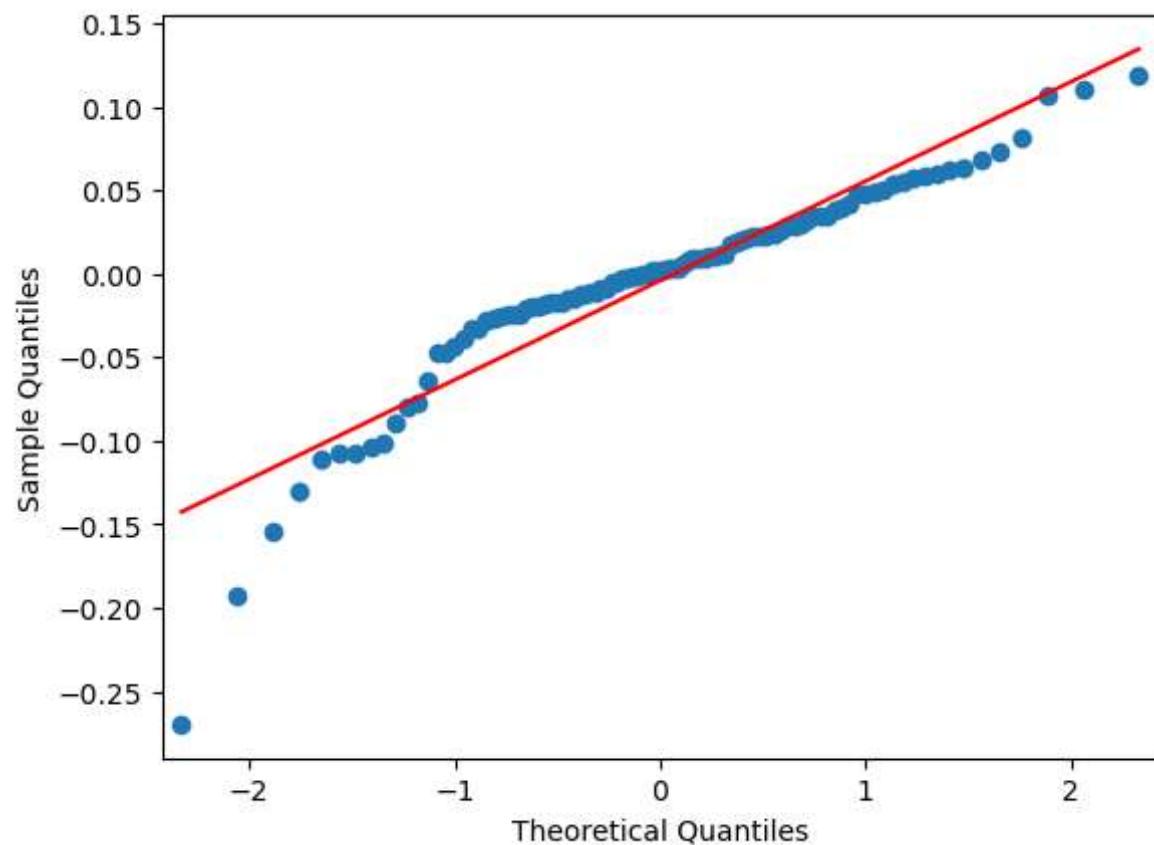
```
In [124]: 1 sns.kdeplot(polyerror)
            2 plt.show()
```



- 1 The observed graph displays a left-skewed distribution, resembling a normal distribution with the mean centered around zero. Notably, there is a cluster of data points on the left side of the graph, suggesting that some students are more inclined to have lower scores affecting their chances of admission.

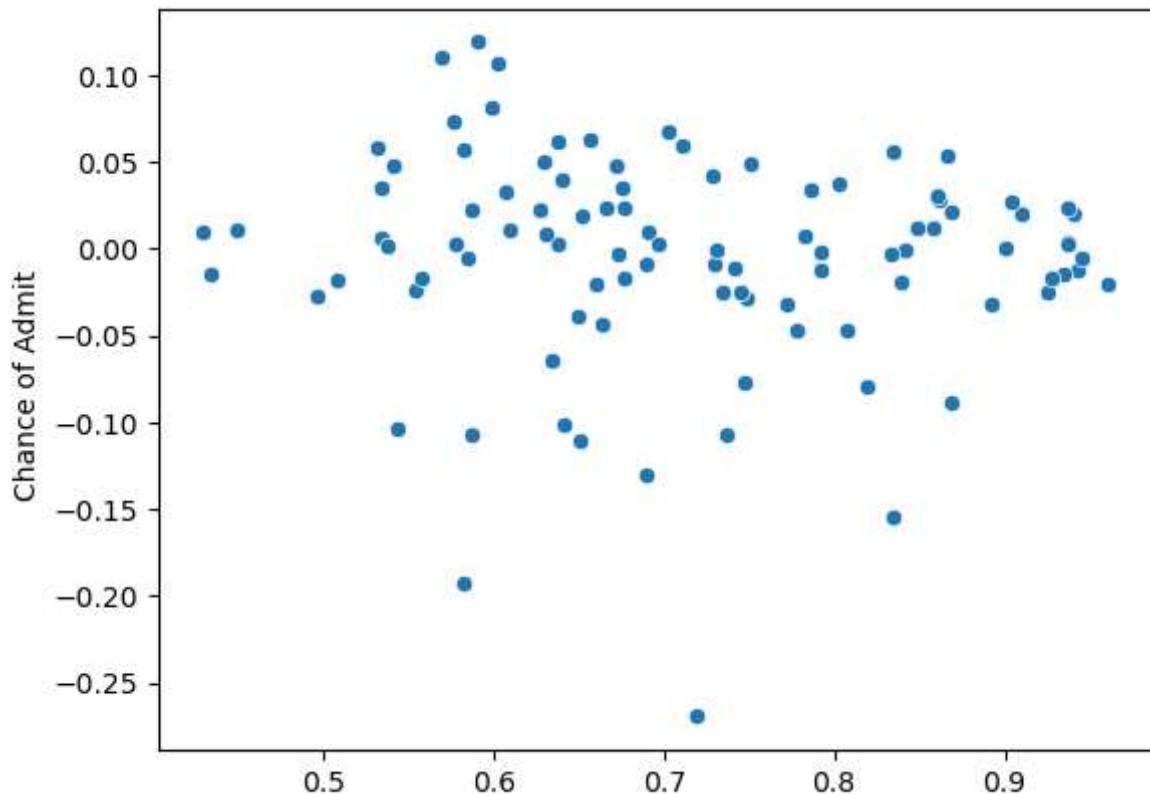
In [125]:

```
1 sm.qqplot(polyerror, line='s')
2 plt.show()
```



Checking for the Homoscedasticity

```
In [126]: 1 sns.scatterplot(x=ypoly_prediction, y=polyerror)
2 plt.show()
```



MSE

```
1 Mean squared error (MSE) serves to quantify the disparity between
actual and predicted values derived from the regression model's best-
fit line. This metric aims to mitigate the discrepancy observed in data
points' distances from the primary regression line. By minimizing this
distance, MSE endeavors to enhance the model's accuracy.
```

```
In [127]: 1 mean_squared_error(y_train,model.predict(Xtrain))
```

Out[127]: 0.00323626117137231

MAE

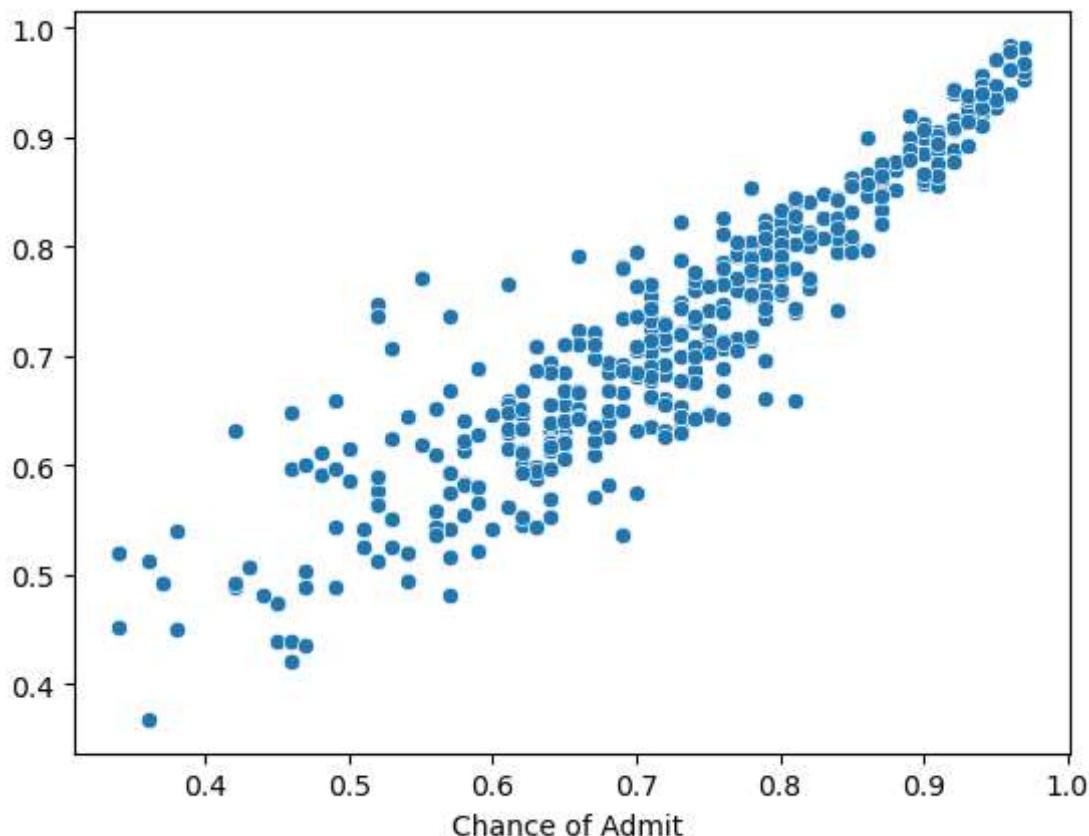
1 Mean squared error (MSE) serves to quantify the disparity between actual and predicted values derived from the regression model's best-fit line. This metric aims to mitigate the discrepancy observed in data points' distances from the primary regression line. By minimizing this distance, MSE endeavors to enhance the model's accuracy.

In [128]: 1 mean_absolute_error(y_train, model.predict(Xtrain))

Out[128]: 0.0400445567153648

RMSE

In [129]: 1 sns.scatterplot(x=y_train, y=model.predict(Xtrain))
2 plt.show()



In [130]: 1 n,d = Xtest.shape
2 np.sqrt(((np.sum(model.predict(Xtest) - y_test)**2)/n))

Out[130]: 0.04070822574379355

```
In [131]: 1 x_sm = sm.add_constant(Xtrain)
2 y_polytrain = y_train.values
3 y_polytrain = y_train1.reshape(-1,1)
4 y_polytrain.shape
```

```
Out[131]: (400, 1)
```

```
In [132]: 1 model = sm.OLS(y_polytrain,x_sm)
2 result = model.fit()
3 print(result.summary())
```

OLS Regression Results

Dep. Variable:	y	R-squared:	0.8			
36						
Model:	OLS	Adj. R-squared:	0.8			
21						
Method:	Least Squares	F-statistic:	54.			
64						
Date:	Tue, 13 Feb 2024	Prob (F-statistic):	5.53e-1			
22						
Time:	14:45:38	Log-Likelihood:	579.			
09						
No. Observations:	400	AIC:	-108			
8.						
Df Residuals:	365	BIC:	-94			
8.5						
Df Model:	34					
Covariance Type:	nonrobust					
5]						
	coef	std err	t	P> t	[0.025	0.97
--						
const	0.3371	0.025	13.344	0.000	0.287	0.3
87						
x1	0.2000	0.091	2.197	0.029	0.021	0.3
79						
x2	0.1231	0.100	1.233	0.218	-0.073	0.3
19						
x3	-0.1228	0.057	-2.169	0.031	-0.234	-0.0
11						
x4	0.0039	0.071	0.054	0.957	-0.136	0.1
44						
x5	0.1283	0.061	2.117	0.035	0.009	0.2
47						
x6	0.4256	0.115	3.715	0.000	0.200	0.6
51						
x7	-0.0131	0.014	-0.919	0.359	-0.041	0.0
15						
x8	0.0443	0.198	0.223	0.823	-0.346	0.4
34						
x9	-0.0661	0.271	-0.244	0.807	-0.598	0.4
66						
x10	0.0609	0.160	0.381	0.704	-0.254	0.3
75						
x11	0.0922	0.194	0.474	0.636	-0.290	0.4
75						
x12	0.1519	0.173	0.879	0.380	-0.188	0.4
92						
x13	-0.4808	0.320	-1.505	0.133	-1.109	0.1
48						
x14	-0.0160	0.075	-0.214	0.831	-0.163	0.1
31						
x15	-0.0378	0.162	-0.234	0.815	-0.356	0.2
80						

x16	0.0032	0.146	0.022	0.983	-0.283	0.2
89						
x17	0.2218	0.187	1.183	0.237	-0.147	0.5
90						
x18	-0.1776	0.133	-1.338	0.182	-0.439	0.0
83						
x19	-0.0221	0.288	-0.077	0.939	-0.588	0.5
44						
x20	0.0429	0.068	0.634	0.526	-0.090	0.1
76						
x21	-0.0023	0.065	-0.035	0.972	-0.131	0.1
26						
x22	0.2839	0.120	2.371	0.018	0.048	0.5
19						
x23	-0.0401	0.108	-0.370	0.712	-0.253	0.1
73						
x24	-0.0838	0.201	-0.417	0.677	-0.479	0.3
11						
x25	0.0073	0.043	0.171	0.865	-0.077	0.0
92						
x26	-0.2593	0.107	-2.424	0.016	-0.470	-0.0
49						
x27	0.0605	0.128	0.473	0.636	-0.191	0.3
12						
x28	-0.0308	0.229	-0.134	0.893	-0.481	0.4
19						
x29	-0.0231	0.050	-0.461	0.645	-0.122	0.0
75						
x30	0.0464	0.070	0.662	0.509	-0.091	0.1
84						
x31	-0.2077	0.192	-1.081	0.281	-0.586	0.1
70						
x32	-0.0093	0.039	-0.236	0.814	-0.087	0.0
68						
x33	0.2789	0.252	1.105	0.270	-0.217	0.7
75						
x34	0.0899	0.079	1.133	0.258	-0.066	0.2
46						
x35	-0.0131	0.014	-0.919	0.359	-0.041	0.0
15						

=====				
==				
Omnibus:	79.274	Durbin-Watson:		2.0
05				
Prob(Omnibus):	0.000	Jarque-Bera (JB):		177.8
52				
Skew:	-1.014	Prob(JB):		2.40e-
39				
Kurtosis:	5.561	Cond. No.		5.23e+
15				
=====				
==				

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
[2] The smallest eigenvalue is 1.26e-28. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.
```

Checking for the Regularisation as per both Lasso and Ridge Regression

- 1 Squared of Weights (regularization term) tries to make all weights to features closer to zero. Since it is not checking on any criteria unlike loss term

Lasso Regression

- 1 Performs L1 regularization, i.e., adds penalty equivalent to the absolute value of the magnitude of coefficients
- 2 Minimization objective = LS Obj + $\alpha * (\text{sum of the absolute value of coefficients})$

```
In [133]: 1 Xtrain.shape, Xtest.shape
```

```
Out[133]: ((400, 35), (100, 35))
```

```
In [134]: 1 y_train.shape, y_test.shape
```

```
Out[134]: ((400,), (100,))
```

```
In [135]: 1 Lassomodel = Lasso(alpha=1)  
2 Lassomodel
```

```
Out[135]:  
└── Lasso  
    └── Lasso(alpha=1)
```

```
In [136]: 1 Lassomodel.fit(Xtrain, y_train)
```

```
Out[136]:  
└── Lasso  
    └── Lasso(alpha=1)
```

```
In [137]: 1 Lassomodel.coef_[:10]
```

```
Out[137]: array([0., 0., 0., 0., 0., 0., 0., 0., 0., 0.])
```

```
In [138]: 1 Lassomodel.intercept_
```

```
Out[138]: 0.7241749999999999
```

```
In [139]: 1 y_test[:10]
```

```
Out[139]: 361    0.93
73     0.84
374    0.39
155    0.77
104    0.74
394    0.89
377    0.47
124    0.57
68     0.68
450    0.82
Name: Chance of Admit , dtype: float64
```

```
In [140]: 1 Lassomodel.predict(Xtest)[:10]
```

```
Out[140]: array([0.724175, 0.724175, 0.724175, 0.724175, 0.724175,
       0.724175, 0.724175, 0.724175, 0.724175])
```

```
In [141]: 1 Lassomodel.score(Xtest, y_test)
```

```
Out[141]: -0.00724844132029312
```

```
In [142]: 1 Lassomodel.score(Xtrain, y_train)
```

```
Out[142]: 0.0
```

```
In [143]: 1 RidgeModelprediction = pd.DataFrame({
2     "Actual" : y_test,
3     "Lasso_predicted": Lassomodel.predict(Xtest)
4 })
5 RidgeModelprediction[:10]
```

```
Out[143]:
```

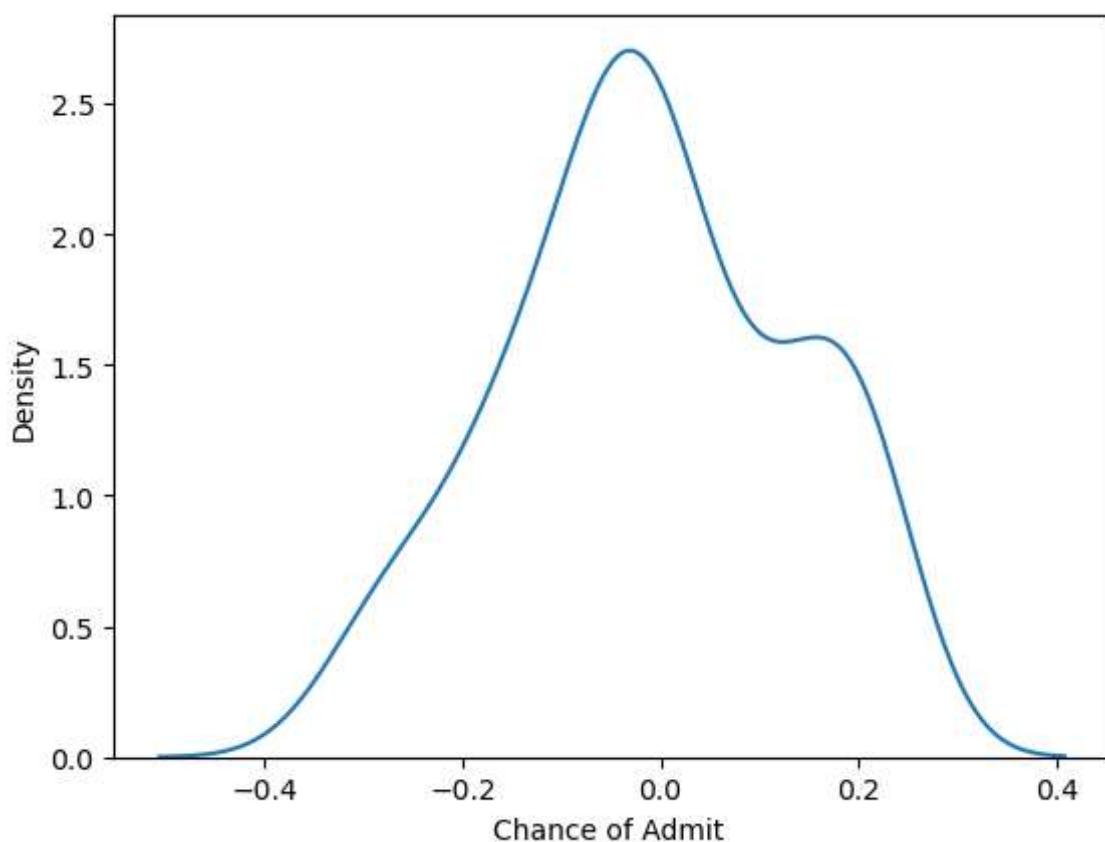
	Actual	Lasso_predicted
361	0.93	0.724175
73	0.84	0.724175
374	0.39	0.724175
155	0.77	0.724175
104	0.74	0.724175
394	0.89	0.724175
377	0.47	0.724175
124	0.57	0.724175
68	0.68	0.724175
450	0.82	0.724175

Residuals

```
In [144]: 1 Lassoerror = y_test - Lassomodel.predict(Xtest)
           2 Lassoerror
```

```
Out[144]: 361    0.205825
           73     0.115825
           374   -0.334175
           155    0.045825
           104    0.015825
           ...
           347   -0.304175
           86    -0.004175
           75    -0.004175
           438   -0.054175
           15    -0.184175
Name: Chance of Admit , Length: 100, dtype: float64
```

```
In [145]: 1 sns.kdeplot(Lassoerror)
           2 plt.show()
```



1 From the Kernel Density Estimation (KDE) graph, it is apparent that the normality of the residuals significantly deteriorates when using the Lasso Regression model. The mean does not align with a consistent point, and the distribution lacks clarity, deviating from the typical Gaussian distribution. This indicates an increase in error and instability in the predictions, leading to less reliable outcomes.

Ridge Regression

- 1 Performs L2 regularization, i.e., adds penalty equivalent to the square of the magnitude of coefficients
- 2 Minimization objective = LS Obj + $\alpha * (\text{sum of square of coefficients})$

In [146]:

```
1 RidgeModel = Ridge(alpha=10)
2 RidgeModel
```

Out[146]:

```
▼ Ridge
Ridge(alpha=10)
```

In [147]:

```
1 RidgeModel.fit(Xtrain, y_train)
```

Out[147]:

```
▼ Ridge
Ridge(alpha=10)
```

In [148]:

```
1 RidgeModel.predict(Xtest).round(3)[:10]
```

Out[148]:

```
array([0.905, 0.786, 0.593, 0.68 , 0.795, 0.858, 0.524, 0.634, 0.795,
       0.795])
```

```
In [149]: 1 RidgeModelprediction = pd.DataFrame({  
2     "Actual" : y_test,  
3     "Ridge_predicted": RidgeModel.predict(Xtest)  
4 })  
5 RidgeModelprediction[:10]
```

Out[149]:

	Actual	Ridge_predicted
361	0.93	0.904978
73	0.84	0.785788
374	0.39	0.592760
155	0.77	0.679749
104	0.74	0.794887
394	0.89	0.857950
377	0.47	0.523715
124	0.57	0.633520
68	0.68	0.795477
450	0.82	0.795446

```
In [150]: 1 RidgeModel.score(Xtrain, y_train)
```

Out[150]: 0.7984222986102354

```
In [151]: 1 RidgeModel.score(Xtest, y_test)
```

Out[151]: 0.808361521766045

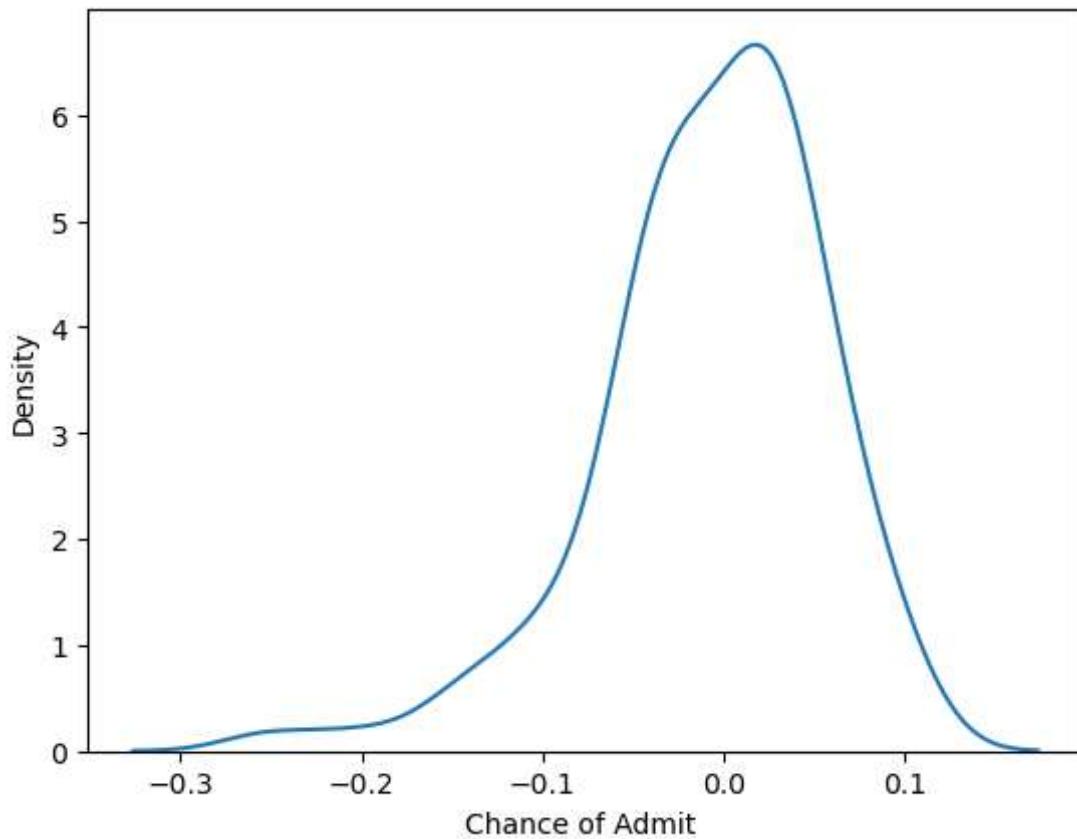
Residuals

```
In [152]: 1 RidgeError = y_test - RidgeModel.predict(Xtest)  
2 RidgeError
```

Out[152]: 361 0.025022
73 0.054212
374 -0.202760
155 0.090251
104 -0.054887
...
347 -0.083784
86 0.037310
75 -0.044379
438 -0.039057
15 -0.102562
Name: Chance of Admit , Length: 100, dtype: float64

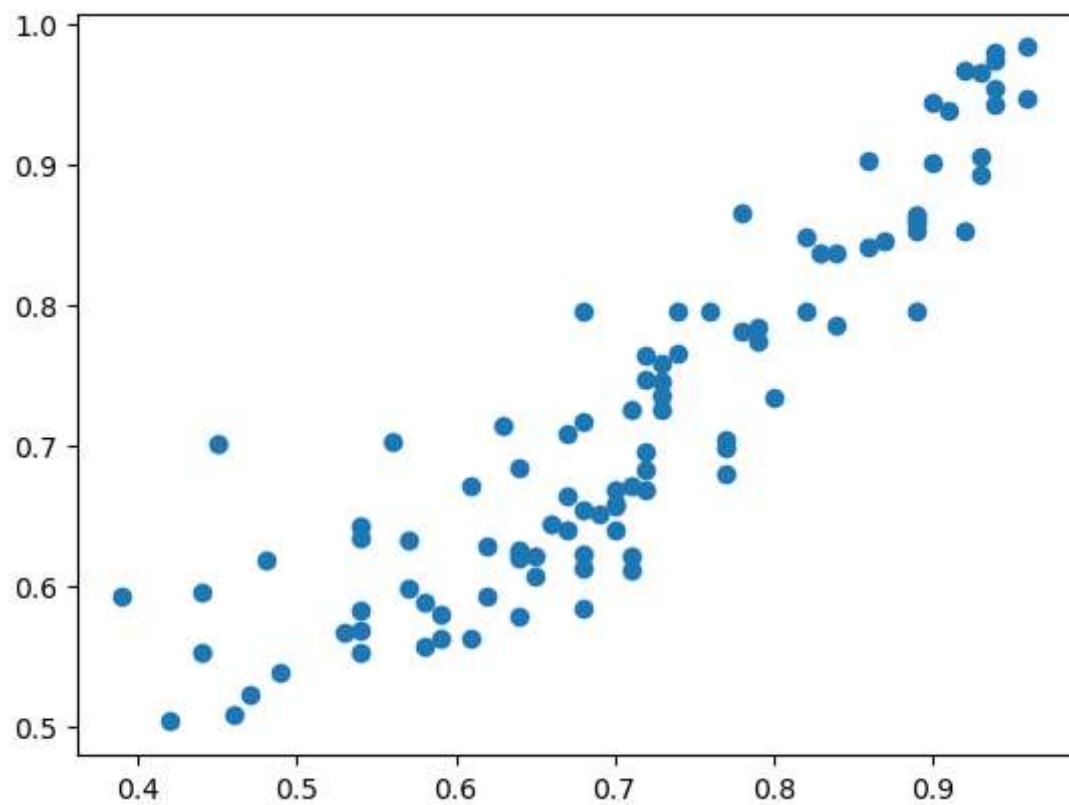
Normality of Residuals

```
In [153]:  
1 sns.kdeplot(RidgeError)  
2 plt.show()
```

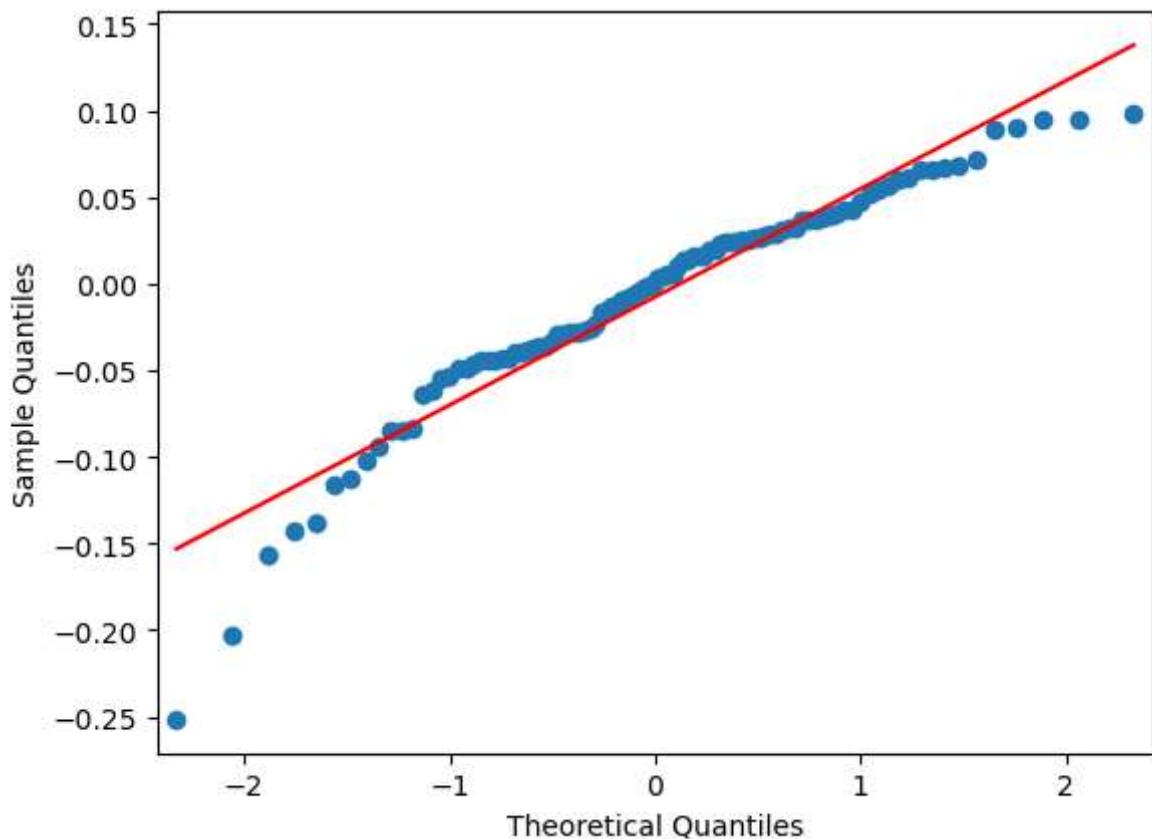


- 1 In the analysis conducted using Ridge regression, it is evident that the scores obtained for both the training and test datasets are slightly lower compared to those achieved using the Linear regression algorithm. Additionally, the Kernel Density Estimation (KDE) plot indicates a slight shift towards the right, suggesting that the mean does not precisely align with zero. This observation implies that the error has a discernible impact, particularly in the context of Ridge regression.

```
In [154]: 1 plt.scatter(x=y_test, y=RidgeModel.predict(Xtest))
2 plt.show()
```



```
In [155]:  
1 sm.qqplot(RidgeError, line="r")  
2 plt.show()
```



- 1 Therefore, our analysis involved utilizing linear regression to assess the data's performance and validity against the underlying assumptions. Additionally, we explored regularization methods such as Lasso and Ridge to potentially enhance the model's predictive accuracy. Despite our efforts, these approaches did not yield satisfactory results in improving performance. Thus, our next step is to employ the Cross Validation Method to further enhance the model's performance.

Cross Validation

Cross Validation to check Performance

```
In [156]: 1 df1 = df  
2 df1.head()
```

Out[156]:

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit	
0	337	118		4	4.5	4.5	9.65	1	0.92
1	324	107		4	4.0	4.5	8.87	1	0.76
2	316	104		3	3.0	3.5	8.00	1	0.72
3	322	110		3	3.5	2.5	8.67	1	0.80
4	314	103		2	2.0	3.0	8.21	0	0.65

```
In [157]: 1 df1.shape
```

Out[157]: (500, 8)

```
In [158]: 1 X.shape, y.shape
```

Out[158]: ((500, 7), (500,))

Cross validation splitting and testing

```
In [159]: 1 X_tr_cv, X_test, y_tr_cv, y_test = train_test_split(X, y, test_size=0.2,  
2 X_train, X_val, y_train, y_val = train_test_split(X_tr_cv, y_tr_cv, test_
```

```
In [160]: 1 print(X_train.shape, y_train.shape)  
2 print(X_val.shape, y_val.shape)  
3 print(X_test.shape, y_test.shape)
```

(300, 7) (300,)
(100, 7) (100,)
(100, 7) (100,)

```
In [161]: 1 ScalingCV = StandardScaler()  
2 ScalingCV
```

Out[161]:

▼ StandardScaler
StandardScaler()

```
In [162]: 1 Xtraincv = ScalingCV.fit_transform(X_train)
           2 Xvalcv = ScalingCV.fit_transform(X_val)
           3 Xtestcv = ScalingCV.transform(X_test)
```

```
In [163]: 1 Modelcv = LinearRegression()
           2 Modelcv
```

Out[163]:

└ LinearRegression

 LinearRegression()

```
In [164]: 1 Modelcv.fit(Xtraincv, y_train)
```

Out[164]:

└ LinearRegression

 LinearRegression()

```
In [165]: 1 Modelcv.coef_
```

Out[165]: array([0.01481262, 0.02647675, 0.00923368, 0.00218789, 0.01131403,
 0.0703194 , 0.00986979])

```
In [166]: 1 Modelcv.intercept_
```

Out[166]: 0.7213333333333334

```
In [167]: 1 Modelcv.score(Xtraincv, y_train)
```

Out[167]: 0.8139677762213995

```
In [168]: 1 y_predcv = Modelcv.predict(Xvalcv)
           2 y_predcv[:10]
```

Out[168]: array([0.86762623, 0.78173248, 0.48973201, 0.77807621, 0.71034476,
 0.75093003, 0.80244193, 0.47546132, 0.90502664, 0.57568781])

```
In [169]: 1 y_val[:10]
```

Out[169]:

235	0.88
105	0.69
367	0.57
76	0.74
251	0.70
3	0.80
48	0.82
457	0.37
46	0.86
493	0.62

Name: Chance of Admit , dtype: float64

```
In [170]: 1 pd.DataFrame({  
2     "Actual":y_val,  
3     "predicted":y_predcv  
4 })[:10]
```

Out[170]:

	Actual	predicted
235	0.88	0.867626
105	0.69	0.781732
367	0.57	0.489732
76	0.74	0.778076
251	0.70	0.710345
3	0.80	0.750930
48	0.82	0.802442
457	0.37	0.475461
46	0.86	0.905027
493	0.62	0.575688

```
In [171]: 1 Modelcv.score(Xvalcv, y_val)
```

Out[171]: 0.8456442461043934

```
In [172]: 1 Modelcv.score(Xtestcv, y_test)
```

Out[172]: 0.8157876639860786

- Utilizing the cross-validation method, we partitioned the dataset into training, testing, and validation subsets to enhance predictive accuracy. Our analysis revealed that the model achieved a performance of 84% on the validation data, showcasing its robustness in capturing the underlying patterns. However, upon evaluating the model on the testing data, the performance slightly declined to 81%, indicating a slight deviation from the linear regression fit. Notably, the Ridge regression approach yielded the most favorable performance, reaching the 84% mark, thus establishing it as the preferred method for optimizing the dataset's performance in the Jamboree project.

Insights

- Prioritize Key Features: Focus on leveraging CGPA, GRE Score, TOEFL Score, and LOR as the top three influential factors for predicting admission chances, with CGPA being the most significant.
-

- 3 Strategic Analysis on CGPA: Conduct an in-depth exploration of the CGPA column to identify universities where students have a high likelihood of admission. Similarly, consider analyzing the GRE and TOEFL Scores to further refine target universities.
- 4
- 5 Instead of binary entries for the Research column, consider quantifying research contributions, such as the number of published papers. While research remains an important factor, it holds less weight compared to CGPA, GRE, and TOEFL Scores.
- 6
- 7 Recognize that University Rating does not significantly impact admission chances. Therefore, focus efforts on other influential factors rather than solely relying on university rankings.

Recommendations

- 1 Feature Importance: Based on the analysis of feature importance, it is evident that both GRE Score and TOEFL Score significantly contribute to predicting the chances of admission. Higher scores in these two features correlate positively with greater chances of admission. Conversely, University Rating appears to have minimal impact on admission predictions. Additionally, the Research feature demonstrates relatively low importance compared to other features in the dataset. Therefore, when considering feature selection for modeling, prioritizing GRE Score and TOEFL Score would likely yield more accurate predictions of admission chances.
- 2
- 3 Model Interpretability : We have employed the Linear Regression algorithm for regression analysis, which yielded the best-fit model with the highest evaluation score for prediction. Our analysis included regularization techniques such as Lasso and Ridge to enhance interpretability and optimize feature selection. However, the linear regression line emerged as the most suitable fit for our dataset.
- 4
- 5 Furthermore, we found that Lasso and Ridge regularization techniques contributed significantly to improving model performance, with a maximum increase of 1% in the evaluation score, resulting in an overall score of 84%. This underscores the importance of these regularization methods in refining feature selection and enhancing predictive accuracy.
- 6
- 7 Preprocessing: We conducted an analysis of the dataset and determined that the feature "Serial No." does not contribute to the predictive power of the model, leading us to eliminate it from consideration. Additionally, we performed thorough outlier treatment, which did not significantly impact the dataset, ensuring data integrity. Standardizing the features resulted in a more uniformly distributed dataset, thereby enhancing the predictive capabilities of our model. These methodical steps underscore the robustness of our approach and demonstrate transparency in our methodology.
- 8

- 9 Model Performance: Among the chosen features, precision emerged as the most reliable metric for model evaluation. Additionally, we analyzed the R2 score and Adjusted R2 score for the remaining features, achieving a consistent range of 83% to 84%. These results signify robust predictive performance and reinforce the effectiveness of our selected features in making accurate predictions.
- 10
- 11 Based on our comprehensive analysis, we have determined that both University Rating and Research have negligible impact on the predictive power of the model. Therefore, we recommend eliminating these features from consideration. Our exploration of various approaches and techniques also indicates that feature extraction is unnecessary in this context. It is evident that University Rating and Research do not contribute significantly to the modeling process and are therefore deemed unsuitable for inclusion in the final model.
- 12
- 13 I conducted an extensive exploratory data analysis (EDA) and feature analysis using various approaches including linear regression, regularization techniques such as Lasso and Ridge, and ElasticNet modeling. Additionally, I experimented with polynomial features to capture non-linear relationships and achieve optimal model fit. Through these methods, I achieved a predictive score of 84%, indicating strong model performance.
- 14
- 15 Furthermore, I employed cross-validation techniques to validate the model and assess its robustness. The resulting R2 score and Adjusted R2 score were consistent, affirming the reliability of the model's predictions. This comprehensive approach ensures accuracy and reliability in our machine learning model for predictive analysis.