# Context

**Jamboree has helped thousands of students make it to top colleges abroad. Be it GMAT, GRE or SAT, their unique problem-solving methods ensure maximum scores with minimum effort.**

**They recently launched a feature where students/learners can come to their website and check their probability of getting into the IVY league college. This feature estimates the chances of graduate admission from an Indian perspective.**

# Problem Statement :

**Help Jamboree in understanding what factors are important in graduate admissions and how these factors are interrelated among themselves. It will also help predict one's chances of admission given the rest of the variables.**

# Column Profiling:

**Serial No. (Unique row ID)**

**GRE Scores (out of 340)**

**TOEFL Scores (out of 120)**

**University Rating (out of 5)**

**Statement of Purpose and Letter of Recommendation Strength (out of 5)**

**Undergraduate GPA (out of 10)**

**Research Experience (either 0 or 1)**

**Chance of Admit (ranging from 0 to 1)**

## Exploratory Data Analysis

## Linear Regression

In [67]:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
from matplotlib import figure
import warnings
warnings.filterwarnings("ignore")
import statsmodels.api as sm
```

In [68]:

```python
data = pd.read_csv("Admission_Predict_Ver1.1.csv")
data.head()
```

Out[68]:

| | 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 [69]:

```
1  data.sample(5)
```

Out[69]:

| | Serial No. | GRE Score | TOEFL Score | University Rating | SOP | LOR | CGPA | Research | Chance of Admit |
|---|---|---|---|---|---|---|---|---|---|
| **364** | 365 | 313 | 102 | 3 | 3.5 | 4.0 | 8.90 | 1 | 0.77 |
| **13** | 14 | 307 | 109 | 3 | 4.0 | 3.0 | 8.00 | 1 | 0.62 |
| **265** | 266 | 313 | 102 | 3 | 2.5 | 2.5 | 8.68 | 0 | 0.71 |
| **236** | 237 | 325 | 112 | 4 | 4.0 | 4.5 | 9.17 | 1 | 0.85 |
| **499** | 500 | 327 | 113 | 4 | 4.5 | 4.5 | 9.04 | 0 | 0.84 |

In [70]:

```
1  # Checking the shape of data
2  data.shape
```

Out[70]:

```
(500, 9)
```

In [71]:

```
1  df = data.copy()
```

In [ ]:

```
1  # Dropping the column which is not required " serial No"
```

In [72]:

```
1  df.drop(["Serial No."],axis = 1 , inplace = True)
```

In [73]:

```
1  # Check null values
2  df.isna().sum()
```

Out[73]:

```
GRE Score            0
TOEFL Score          0
University Rating    0
SOP                  0
LOR                  0
CGPA                 0
Research             0
Chance of Admit      0
dtype: int64
```

In [74]:

```
1  # Information about the data type of all columns
2  df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 8 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   GRE Score          500 non-null    int64
 1   TOEFL Score        500 non-null    int64
 2   University Rating  500 non-null    int64
 3   SOP                500 non-null    float64
 4   LOR                500 non-null    float64
 5   CGPA               500 non-null    float64
 6   Research           500 non-null    int64
 7   Chance of Admit    500 non-null    float64
dtypes: float64(4), int64(4)
memory usage: 31.4 KB
```

**There is no null values in any column.**

**No null value detected.**

In [75]:

```
1  # checking for the number of unique values in each columns
2  df.nunique()
```

Out[75]:

```
GRE Score            49
TOEFL Score          29
University Rating     5
SOP                   9
LOR                   9
CGPA                184
Research              2
Chance of Admit      61
dtype: int64
```

## Observation

**University Rating,SOP,LOR,Research are seems to be categorical variables as the number of unique values are very small.**

**rest of the features are numeric , and ordinal . (University Rating,SOP,LOR,Research are discrete ) and rest are continuous**

**also if SOP , University rating , LOR and research can be considered as numeric ordinal data.**
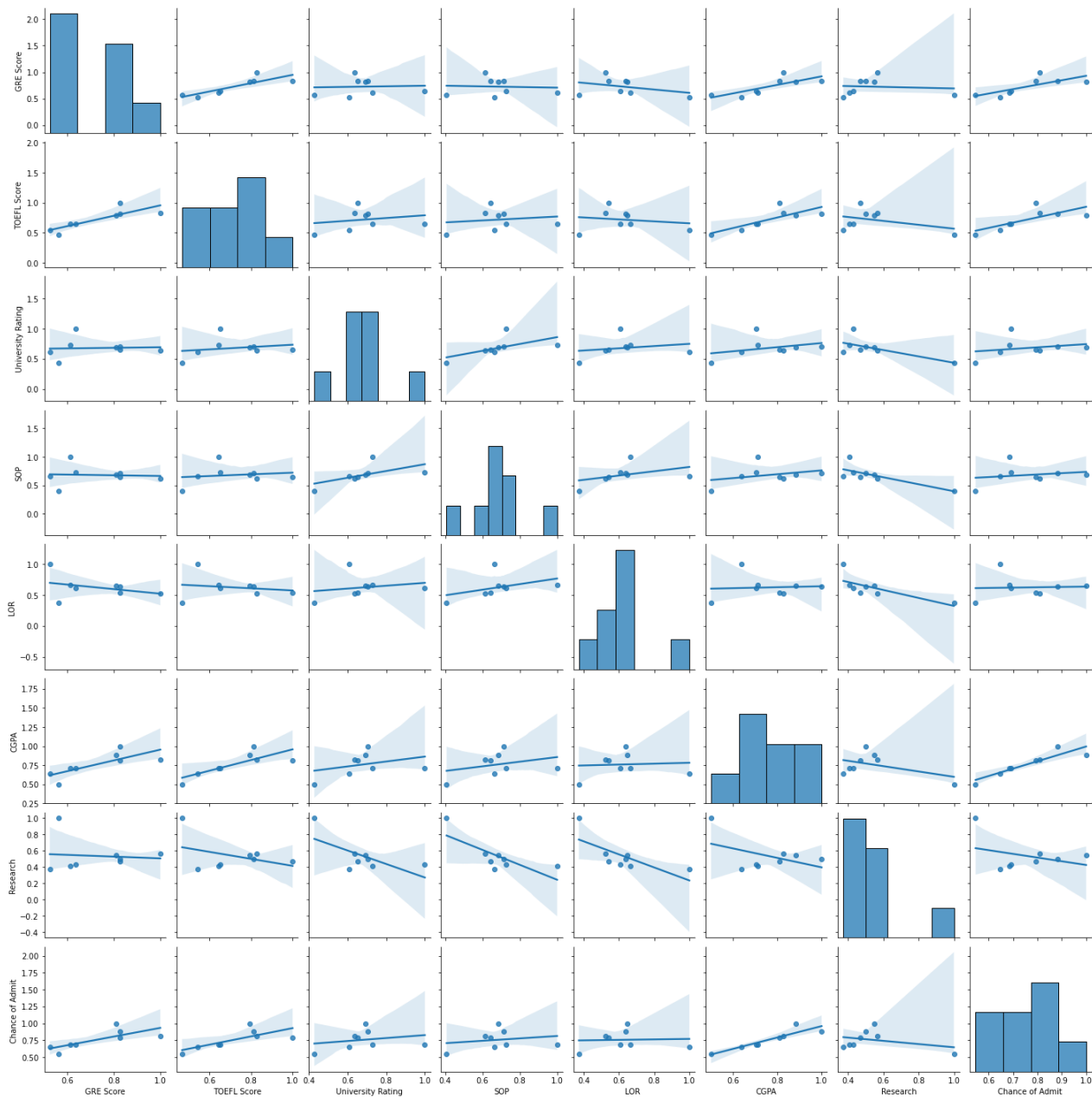
In [ ]:

```
1
```

**Checking the overall linearity and correlation across all features using pairplot**

In [76]:

```
1  sns.pairplot(df.corr(),kind = "reg" )
```

Out[76]:

`<seaborn.axisgrid.PairGrid at 0x217e6899940>`
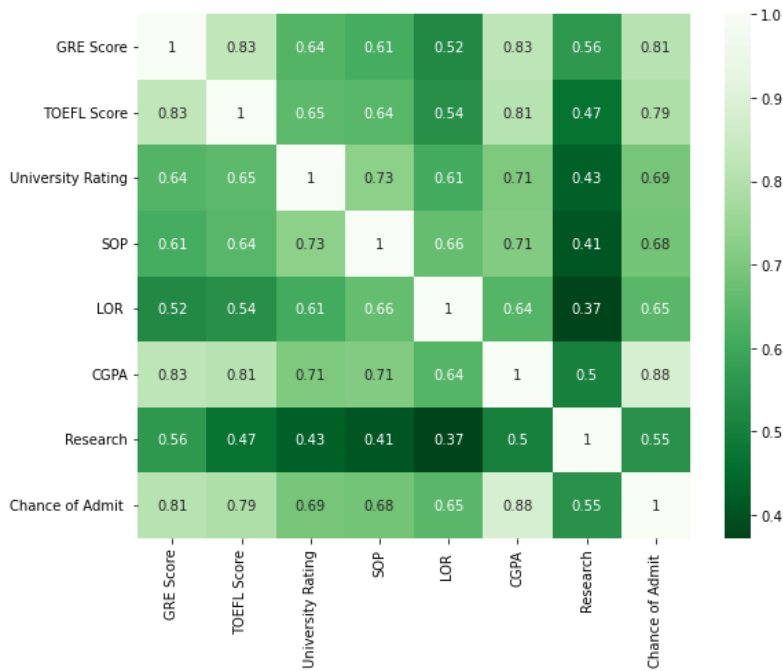


# Overall look at the Correlation :

In [77]:

```
1
2 plt.figure(figsize=(9,7))
3 sns.heatmap(df.corr(),annot=True,cmap = "Greens_r")
```

Out[77]:

`<AxesSubplot:>`



## Observation

**Independent Variables (Input data): GRE Score, TOEFL Score, University Rating, SOP, LOR, CGPA, Research**

**Target/Dependent Variable : Chance of Admit (the value we want to predict)**

**from above correlation heatmap , we can observe GRE score TOEFL score and CGPA have very high correlation with Change of admission.**

**University rating, SOP ,LOR and Research have comparatively slightly less correlated than other features.**

In [ ]:

```
1
```

In [78]:

```
1 # For better convenience remove the spaces from the column names
2 df.columns
```

Out[78]:

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

In [79]:

```
1 df.columns = ['GRE_Score', 'TOEFL_Score', 'University_Rating', 'SOP', 'LOR', 'CGPA',
2         'Research', 'Chance_of_Admit']
```

In [80]:

```
1  df.sample(5)
```

Out[80]:

|  | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Chance_of_Admit |
|---|---|---|---|---|---|---|---|---|
| 463 | 304 | 107 | 3 | 3.5 | 3.0 | 7.86 | 0 | 0.57 |
| 158 | 306 | 106 | 2 | 2.0 | 2.5 | 8.14 | 0 | 0.61 |
| 110 | 305 | 108 | 5 | 3.0 | 3.0 | 8.48 | 0 | 0.61 |
| 288 | 314 | 104 | 4 | 5.0 | 5.0 | 9.02 | 0 | 0.82 |
| 375 | 304 | 101 | 2 | 2.0 | 2.5 | 7.66 | 0 | 0.38 |

In [ ]:

```
1
```

# Outliers in the Data

In [81]:

```
1  def detect_outliers(data):
2      length_before = len(data)
3      Q1 = np.percentile(data,25)
4      Q3 = np.percentile(data,75)
5      IQR = Q3-Q1
6      upperbound = Q3+1.5*IQR
7      lowerbound = Q1-1.5*IQR
8      if lowerbound < 0:
9          lowerbound = 0
10
11     length_after = len(data[(data>lowerbound)&(data<upperbound)])
12     return f"{np.round((length_before-length_after)/length_before,4)} % Outliers data from input data found"
13
14
```

In [82]:

```
1  for col in df.columns:
2      print(col," : ",detect_outliers(df[col]))
```

```
GRE_Score  :  0.0 % Outliers data from input data found
TOEFL_Score  :  0.0 % Outliers data from input data found
University_Rating  :  0.0 % Outliers data from input data found
SOP  :  0.0 % Outliers data from input data found
LOR  :  0.024 % Outliers data from input data found
CGPA  :  0.0 % Outliers data from input data found
Research  :  0.44 % Outliers data from input data found
Chance_of_Admit  :  0.004 % Outliers data from input data found
```

In [83]:

```
1  detect_outliers(df)
```

Out[83]:

```
'0.0 % Outliers data from input data found'
```

# Observation

**there are no significant amount of outliers found in the data.**

In [ ]:

```
1
```

# Descriptive analysis of all numerical features :

In [84]:

```
1 df.describe()
```

Out[84]:

|  | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Chance_of_Admit |
|---|---|---|---|---|---|---|---|---|
| count | 500.000000 | 500.000000 | 500.000000 | 500.000000 | 500.000000 | 500.000000 | 500.000000 | 500.00000 |
| mean | 316.472000 | 107.192000 | 3.114000 | 3.374000 | 3.48400 | 8.576440 | 0.560000 | 0.72174 |
| std | 11.295148 | 6.081868 | 1.143512 | 0.991004 | 0.92545 | 0.604813 | 0.496884 | 0.14114 |
| min | 290.000000 | 92.000000 | 1.000000 | 1.000000 | 1.00000 | 6.800000 | 0.000000 | 0.34000 |
| 25% | 308.000000 | 103.000000 | 2.000000 | 2.500000 | 3.00000 | 8.127500 | 0.000000 | 0.63000 |
| 50% | 317.000000 | 107.000000 | 3.000000 | 3.500000 | 3.50000 | 8.560000 | 1.000000 | 0.72000 |
| 75% | 325.000000 | 112.000000 | 4.000000 | 4.000000 | 4.00000 | 9.040000 | 1.000000 | 0.82000 |
| max | 340.000000 | 120.000000 | 5.000000 | 5.000000 | 5.00000 | 9.920000 | 1.000000 | 0.97000 |

# Observation :

**chances of admit is a probability measure , which is within 0 to 1 which is good (no outliers or missleading data in column).**

**Range of GRE score looks like between 290 to 340.**

**range of TOEFL score is between 92 to 120.**

**university rating , SOP and LOR are distributed between range of 1 to 5.**

**CGPA range is between 6.8 to 9.92.**

In [ ]:

```
1
```

Type *Markdown* and LaTeX: $\alpha^2$

# Graphical Analysis :
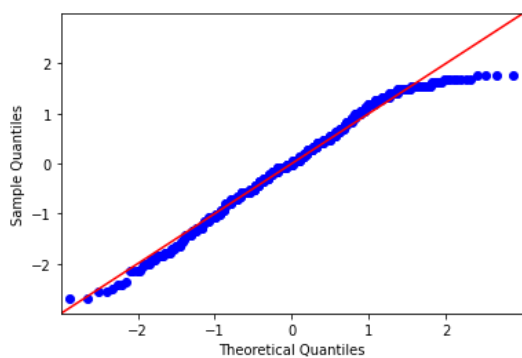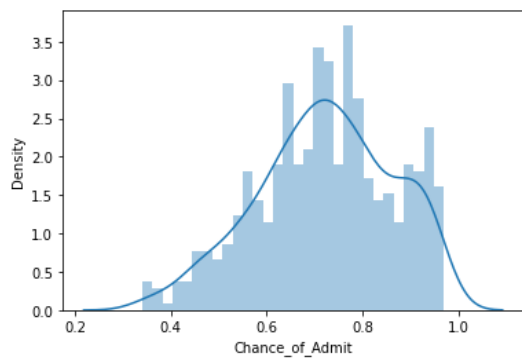
## Distributions / Histogram and count plot :

In [ ]:

```
1
```

## Chance_of_Admit

In [17]:

```python
sns.distplot(df["Chance_of_Admit"],bins = 30)
sm.qqplot(df["Chance_of_Admit"],fit=True, line="45")
plt.show()
```
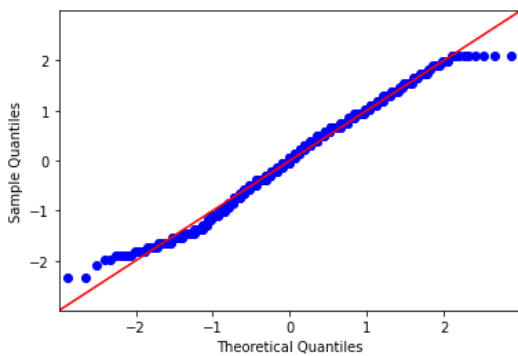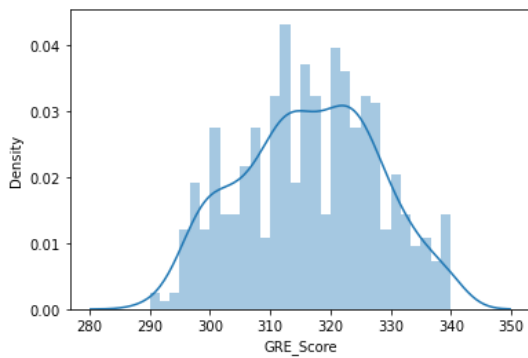




In [ ]:

```python

```

**GRE_Score**

In [18]:

```
1
2  sns.distplot(df["GRE_Score"], bins = 30)
3  sm.qqplot(df["GRE_Score"],fit=True, line="45")
4  plt.show()
```
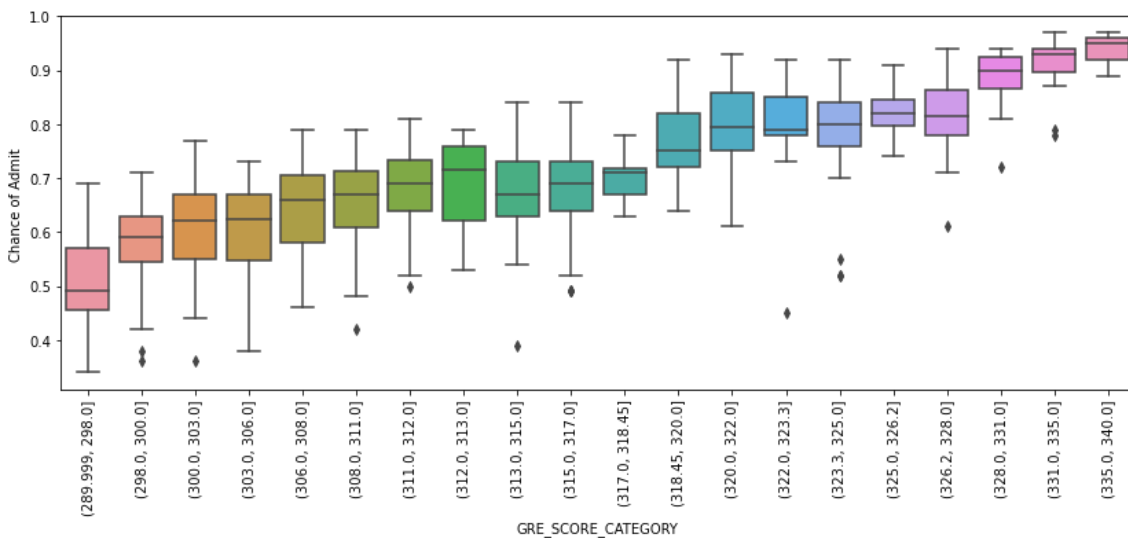




In [ ]:

```
1
```

In [19]:

```
1  data["GRE_SCORE_CATEGORY"]=pd.qcut(data["GRE Score"],20)
2  plt.figure(figsize=(14,5))
3  sns.boxplot(y = data["Chance of Admit "], x = data["GRE_SCORE_CATEGORY"])
4  plt.xticks(rotation = 90)
5  plt.show()
6
```



# Observation :

**From above boxplot (distribution of chance of admition (probability of getting admition) as per GRE score ) :**
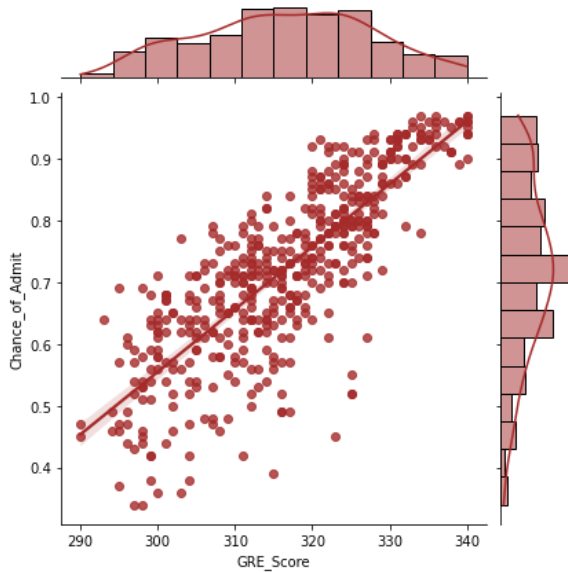
**with higher GRE score , there is high probability of getting an admition .**

In [20]:

```
1  sns.jointplot(df["GRE_Score"],df["Chance_of_Admit"], kind = "reg",color = "brown" )
2
```

Out[20]:

```
<seaborn.axisgrid.JointGrid at 0x217ecfd5df0>
```
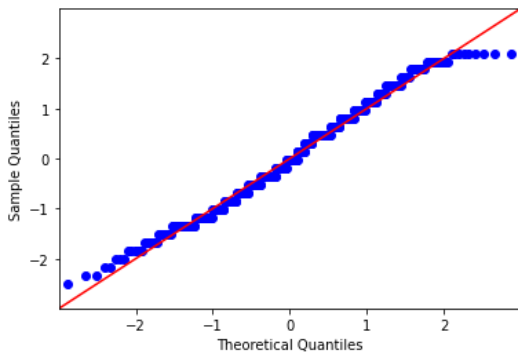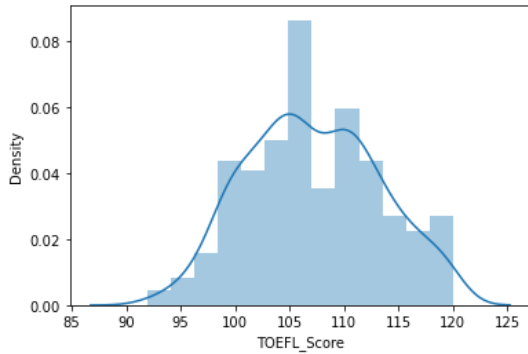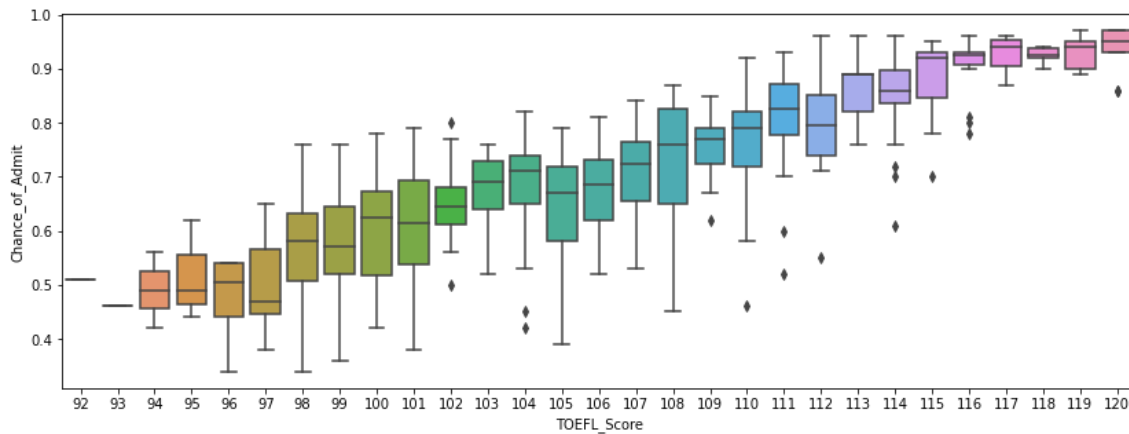
**TOEFL_Score**

In [21]:

```
1  # TOEFL_Score
2
3  sns.distplot(df["TOEFL_Score"])
4  sm.qqplot(df["TOEFL_Score"],fit=True, line="45")
5  plt.show()
6  plt.figure(figsize=(14,5))
7  sns.boxplot(y = df["Chance_of_Admit"], x = df["TOEFL_Score"])
8
```





Out[21]:

```
<AxesSubplot:xlabel='TOEFL_Score', ylabel='Chance_of_Admit'>
```



# Observation :

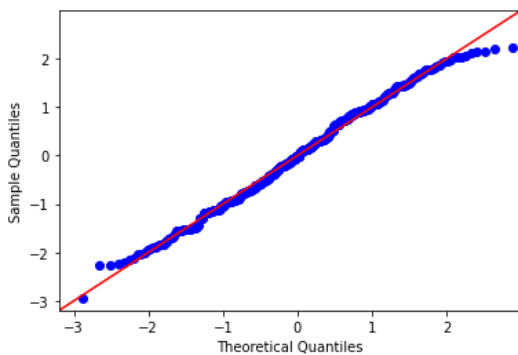**Students having high toefl score , has higher probability of getting admition .**

In [ ]:

```
1
```

In [ ]:

```
1
```

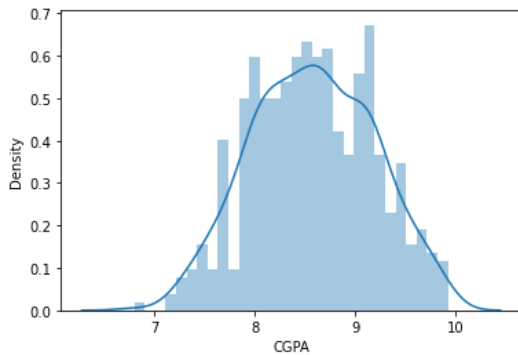**CGPA**

In [22]:

```
1
2
3  sns.distplot(df["CGPA"], bins = 30)
4  sm.qqplot(df["CGPA"],fit=True, line="45")
5  plt.show()
6
```





# Observation :

**Chance of admit and GRE score are nearly normally distrubted.**

**GRE score, TOEFL score and CGPA has a strong correlation with chance of addmission .**

In [ ]:

```
1
```

In [ ]:

```
1
```

# Distribution of all other categorical features :

In [23]:

```python
plt.figure(figsize=(15,10))
plt.subplot(2,2,1)
sns.countplot(df["University_Rating"])
plt.subplot(2,2,2)
sns.countplot(df["LOR"])
plt.subplot(2,2,3)
sns.countplot(df["SOP"])
plt.subplot(2,2,4)
sns.countplot(df["Research"])
```

Out[23]:

```
<AxesSubplot:xlabel='Research', ylabel='count'>
```

In [ ]:

```
1
```

In [24]:

```python
sns.pairplot(df,y_vars = ["Chance_of_Admit"])
plt.title("Pair plot Chance of admit vs all the features")
plt.show()
```

In [ ]:

```
1
```

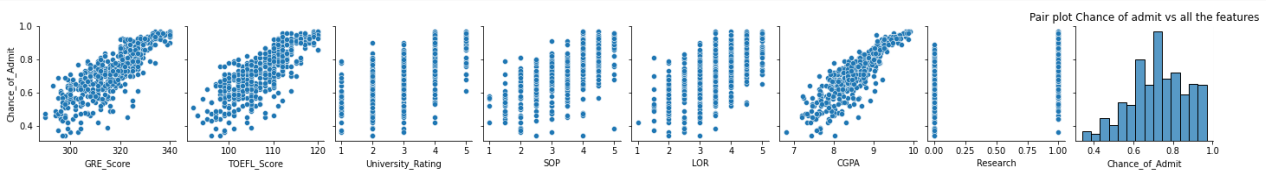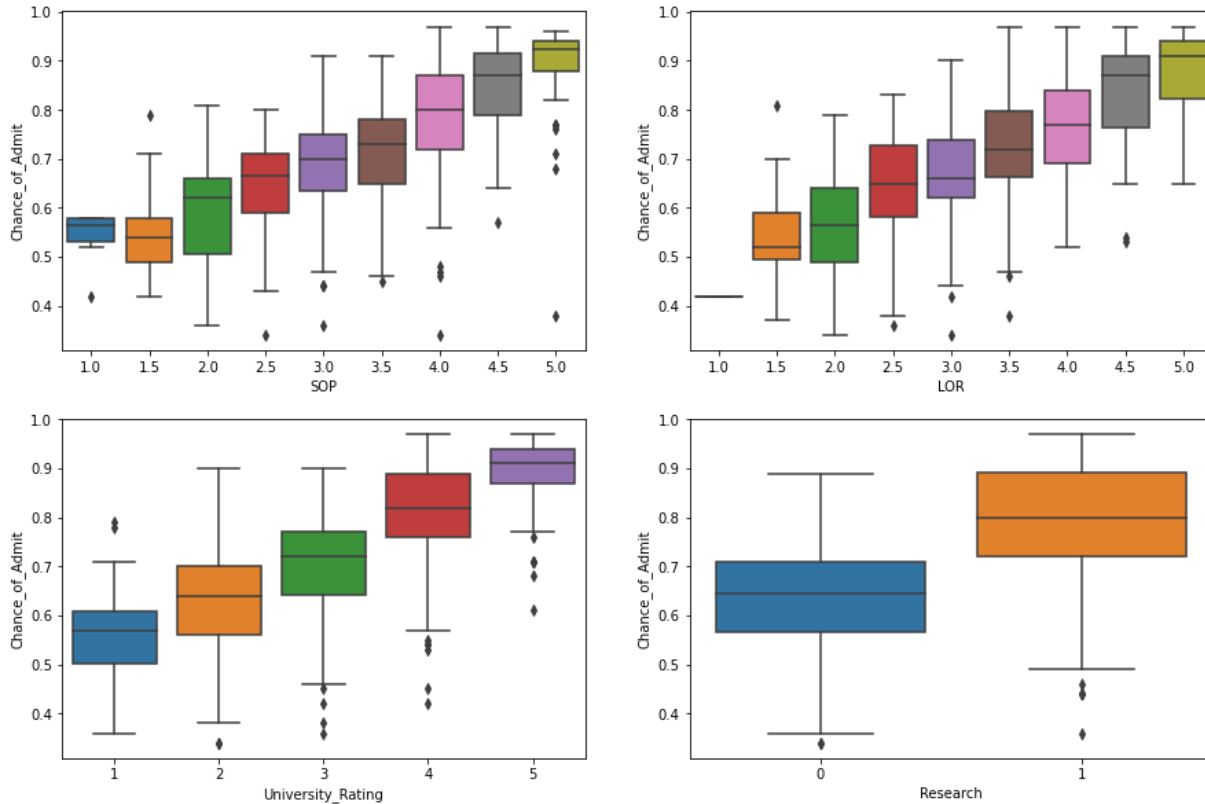# Categorical features - vs - chances of admission boxplot :

In [25]:

```python
plt.figure(figsize=(15,10))
plt.subplot(2,2,1)
sns.boxplot(y = df["Chance_of_Admit"], x = df["SOP"])
plt.subplot(2,2,2)
sns.boxplot(y = df["Chance_of_Admit"], x = df["LOR"])
plt.subplot(2,2,3)
sns.boxplot(y = df["Chance_of_Admit"], x = df["University_Rating"])
plt.subplot(2,2,4)
sns.boxplot(y = df["Chance_of_Admit"], x = df["Research"])
plt.show()
```

In [ ]:

```
1
```

## Observation:

from above plots, we can observe , statement of purpose SOP strength is positively correlated with Chance of Admission .

we can also similar pattern in Letter of Recommendation Stength and University rating , have positive correlation with Chaces of Admission .

Student having research has higher chances of Admission , but also we can observe some outliers within that caregory.

In [ ]:

```
1
```

In [ ]:

```
1
```

## Linearity : How features are correlated with Target variable - chance of admit :

In [26]:

```
for col in df.columns[:-1]:
    print(col)
    plt.figure(figsize=(3,3))
    sns.jointplot(df[col],df["Chance_of_Admit"],kind="reg"  )
    plt.show()
```

GRE_Score

<Figure size 216x216 with 0 Axes>



TOEFL_Score

<Figure size 216x216 with 0 Axes>



University_Rating

<Figure size 216x216 with 0 Axes>

SOP

```
<Figure size 216x216 with 0 Axes>
```



LOR

```
<Figure size 216x216 with 0 Axes>
```



CGPA

```
<Figure size 216x216 with 0 Axes>
```

Research

```
<Figure size 216x216 with 0 Axes>
```



**In [ ]:**
```
1
```

**Line**

**In [ ]:**
```
1
```

**In [15]:**
```
1  from sklearn.preprocessing import StandardScaler
2
3  from sklearn.linear_model import LinearRegression
4  from sklearn.model_selection import train_test_split
5
6  from statsmodels.stats.outliers_influence import variance_inflation_factor
7
8  from sklearn.metrics import r2_score,mean_squared_error,mean_absolute_error, adjusted_mutual_info_score
9  from sklearn.feature_selection import f_regression
```

**In [ ]:**
```
1
```

**In [ ]:**
```
1
```

**In [ ]:**
```
1
```

**In [152]:**
```
1
2  X = df.drop(["Chance_of_Admit"],axis = 1)  # independent variables
3  y = df["Chance_of_Admit"].values.reshape(-1,1) # target / dependent variables
4
5
```

**In [ ]:**
```
1
```

## Standardising data

In [153]:

```
1  standardizer = StandardScaler()
2  standardizer.fit(X)
3  x = standardizer.transform(X)   # standardising the data
4
```

## test train spliting :

In [ ]:

```
1
```

In [154]:

```
1  X_train , X_test, y_train , y_test = train_test_split(x,y,random_state = 1,test_size = 0.2 )   # test train split
2
3
```

In [155]:

```
1  # after spliting, checking for the shape of test and  train data
2
3
4  X_train.shape,X_test.shape
```

Out[155]:

```
((400, 7), (100, 7))
```

In [156]:

```
1  y_train.shape, y_test.shape
2
```

Out[156]:

```
((400, 1), (100, 1))
```

# training the model

In [157]:

```
1  LinearRegression = LinearRegression()      # training LinearRegression model
2  LinearRegression.fit(X_train,y_train)
3
```

Out[157]:

```
LinearRegression()
```

## R2 score on train data :

In [158]:

```
1  r2_score(y_train,LinearRegression.predict(X_train))
```

Out[158]:

```
0.8215099192361265
```

## R2 score on test data :

In [159]:

```
1  r2_score(y_test,LinearRegression.predict(X_test) )
2
```

Out[159]:

```
0.8208741703103732
```

```
In [ ]:
1
```

# All the feature's coefficients and Intercept :

```
In [160]:
1 ws = pd.DataFrame(LinearRegression.coef_.reshape(1,-1),columns=df.columns[:-1])
2 ws["Intercept"] = LinearRegression.intercept_
3 ws
4
```

Out[160]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020675 | 0.019284 | 0.007001 | 0.002975 | 0.013338 | 0.070514 | 0.009873 | 0.722881 |

```
In [161]:
1 LinearRegression_Model_coefs = ws
2 LinearRegression_Model_coefs
3
```

Out[161]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020675 | 0.019284 | 0.007001 | 0.002975 | 0.013338 | 0.070514 | 0.009873 | 0.722881 |

```
In [162]:
1
2 def AdjustedR2score(R2,n,d):
3     return 1-(((1-R2)*(n-1))/(n-d-1))
4
```

```
In [163]:
1
2 y_pred = LinearRegression.predict(X_test)
3
4 print("MSE:",mean_squared_error(y_test,y_pred)) # MSE
5 print("RMSE:",np.sqrt(mean_squared_error(y_test,y_pred))) #RMSE
6 print("MAE :",mean_absolute_error(y_test,y_pred) ) # MAE
7 print("r2_score:",r2_score(y_test,y_pred)) # r2score
8 print("Adjusted R2 score :", AdjustedR2score(r2_score(y_test,y_pred),len(X),X.shape[1]))  # adjusted R2 score
9
10
```

```
MSE: 0.0034590988971363833
RMSE: 0.058814104576507695
MAE : 0.040200193804157944
r2_score: 0.8208741703103732
Adjusted R2 score : 0.8183256320830818
```

# Using Sklearn | Stochastic Gradient Descent Aalgorithm"

```
In [98]:
1 X = df.drop(["Chance_of_Admit"],axis = 1)
2 y = df["Chance_of_Admit"]
3 X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=2)
```

```
In [99]:
1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3
```

```
In [100]:
1 scaler.fit(X_train)
```

Out[100]:

```
StandardScaler()
```

In [101]:

```
1  X_train = scaler.transform(X_train)
2  X_test = scaler.transform(X_test) # apply same transformation to test data
3
```

In [102]:

```
1  from sklearn.linear_model import SGDRegressor
2  from sklearn.pipeline import make_pipeline
3  sgd = make_pipeline(StandardScaler(), SGDRegressor(max_iter=1000, tol=1e-3))
```

In [103]:

```
1  sgd.fit(X_train, y_train)
2
```

Out[103]:

```
Pipeline(steps=[('standardscaler', StandardScaler()),
                ('sgdregressor', SGDRegressor())])
```

In [104]:

```
1  y_pred = sgd.predict(X_test)
2
```

In [105]:

```
1  y_test = y_test.values
```

In [106]:

```
1  r2_score(y_test,y_pred)
```

Out[106]:

```
0.782989764614124
```

In [247]:

```
1
2  # trying different algorithms and different variations with features.
```

# Linear Regression using Statsmodel library

In [115]:

```
1  import statsmodels.api as sm
2  X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=2)
3
4
5
```

In [116]:

```
1  X_train_sm = X_train
2  X_test_sm = X_test
```

In [117]:

```
1  X_train_sm = sm.add_constant(X_train_sm)
2  X_test_sm = sm.add_constant(X_test_sm)
```

# Assumptions of linear regression

## 1)No multicollinearity

## 2)The mean of residual is nearly zero.

## 3)Linearity of Variables

## 4)Test of homoscedasticity

## 5)Normality of residual

In [ ]:

```
1
```

## *Multicollinearity check :

### checking vif scores :

In [49]:

```
1  vifs = []
2
3  for i in range(X_train.shape[1]):
4
5      vifs.append((variance_inflation_factor(exog = X_train,
6                                    exog_idx=i)))
7  vifs
```

Out[49]:

```
[4.244635042406759,
 4.063329028077076,
 2.5932198746362025,
 2.7053574355882417,
 1.9762313259255433,
 4.766976206732742,
 1.466566895741921]
```

In [51]:

```
1  pd.DataFrame({ "coef_name : " : X.columns ,
2              "vif : ": np.around(vifs,2)})
```

Out[51]:

|   | coef_name : | vif : |
|---|---|---|
| 0 | GRE_Score | 4.24 |
| 1 | TOEFL_Score | 4.06 |
| 2 | University_Rating | 2.59 |
| 3 | SOP | 2.71 |
| 4 | LOR | 1.98 |
| 5 | CGPA | 4.77 |
| 6 | Research | 1.47 |

In [118]:

```
1  olsres = sm.OLS(y_train,X_train_sm).fit()
```

In [119]:

```
1  print(olsres.summary())
```

```
                          OLS Regression Results
==============================================================================
Dep. Variable:        Chance_of_Admit   R-squared:                       0.829
Model:                            OLS   Adj. R-squared:                  0.826
Method:                 Least Squares   F-statistic:                     272.1
Date:                Wed, 14 Dec 2022   Prob (F-statistic):          3.33e-146
Time:                        18:27:43   Log-Likelihood:                 573.41
No. Observations:                 400   AIC:                            -1131.
Df Residuals:                     392   BIC:                            -1099.
Df Model:                           7
Covariance Type:            nonrobust
=====================================================================================
                        coef    std err          t      P>|t|      [0.025      0.975]
-------------------------------------------------------------------------------------
const                -1.3418      0.116    -11.613      0.000      -1.569      -1.115
GRE_Score             0.0021      0.001      3.893      0.000       0.001       0.003
TOEFL_Score           0.0030      0.001      3.024      0.003       0.001       0.005
University_Rating     0.0048      0.004      1.185      0.237      -0.003       0.013
SOP                   0.0021      0.005      0.428      0.669      -0.008       0.012
LOR                   0.0186      0.005      4.131      0.000       0.010       0.027
CGPA                  0.1134      0.011     10.633      0.000       0.092       0.134
Research              0.0247      0.007      3.476      0.001       0.011       0.039
==============================================================================
Omnibus:                       94.166   Durbin-Watson:                   1.943
Prob(Omnibus):                  0.000   Jarque-Bera (JB):              231.309
Skew:                          -1.158   Prob(JB):                     5.92e-51
Kurtosis:                       5.918   Cond. No.                     1.33e+04
==============================================================================

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.33e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

In [120]:

```
1  r2_score(y_test,olsres.predict(X_test_sm))
```

Out[120]:

0.7927524897595928

In [ ]:

```
1
```

# Observation :

## VIF score are all below 5 , doesnt seem to have very high multicolinearity.

## same result of r2 value , as sklearn OLS regressor. ,

In [ ]:

```
1
```

In [ ]:

```
1
```

## Residual analysis :

In [164]:

```
1  y_predicted = LinearRegression.predict(X_train)
2  y_predicted.shape
3
4
```
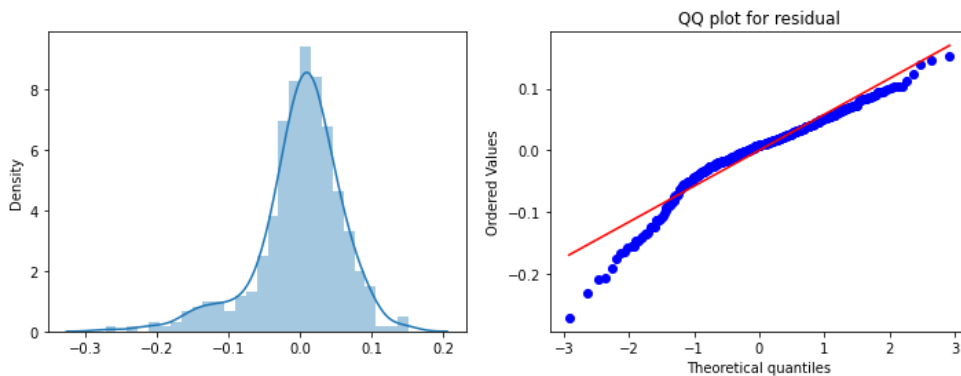
Out[164]:

(400, 1)

In [165]:

```
1  residuals = (y_train - y_predicted)
2  plt.figure(figsize=(12,4))
3  plt.subplot(1,2,1)
4  sns.distplot(residuals)
5  plt.subplot(1,2,2)
6  stats.probplot(residuals.reshape(-1,), plot = plt)
7  plt.title('QQ plot for residual')
8  plt.show()
```



In [ ]:

```
1
```

## Linearity of varibales

In [166]:

```
1  sns.pairplot(df,y_vars = ["Chance_of_Admit"])
2  plt.show()
```



In [ ]:

```
1
```

## Test of homoscedasticity | plotting y_predicted and residuals

In [173]:

```
1  # Test of homoscedasticity
2  sns.scatterplot(y_predicted.reshape(-1,), residuals.reshape(-1,))
3  plt.xlabel('y_predicted')
4  plt.ylabel('Residuals')
5  plt.axhline(y=0)
6  plt.title("Y_predicted vs residuals, check of homoscedasticity")
7  plt.show()
```

In [ ]:

```
1
```

# Observation

## Homoscedasticity

**from above residual plot , we can observe the varinace is not so constant .**

**all residuals are not evenly distributed.**

In [ ]:

```
1
```

# Model Regularisation :

In [175]:

```
1  from sklearn.linear_model import Ridge   # L2 regualrization
2  from sklearn.linear_model import Lasso   # L1 regualrization
3  from sklearn.linear_model import ElasticNet
```
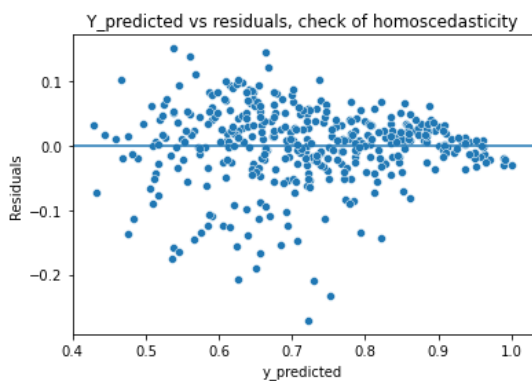
In [ ]:

```
1
```

# L2 regularization

## Ridge regression :

In [178]:

```
1  ## Hyperparameter Tuning : for appropriate lambda value :
2
3  train_R2_score = []
4  test_R2_score = []
5  lambdas = []
6  train_test_difference_Of_R2 =  []
7  lambda_ = 0
8  while lambda_ <= 5:
9      lambdas.append(lambda_)
10     RidgeModel = Ridge(lambda_)
11     RidgeModel.fit(X_train,y_train)
12     trainR2 = RidgeModel.score(X_train,y_train)
13     testR2 = RidgeModel.score(X_test,y_test)
14     train_R2_score.append(trainR2)
15     test_R2_score.append(testR2)
16
17     lambda_ += 0.01
```
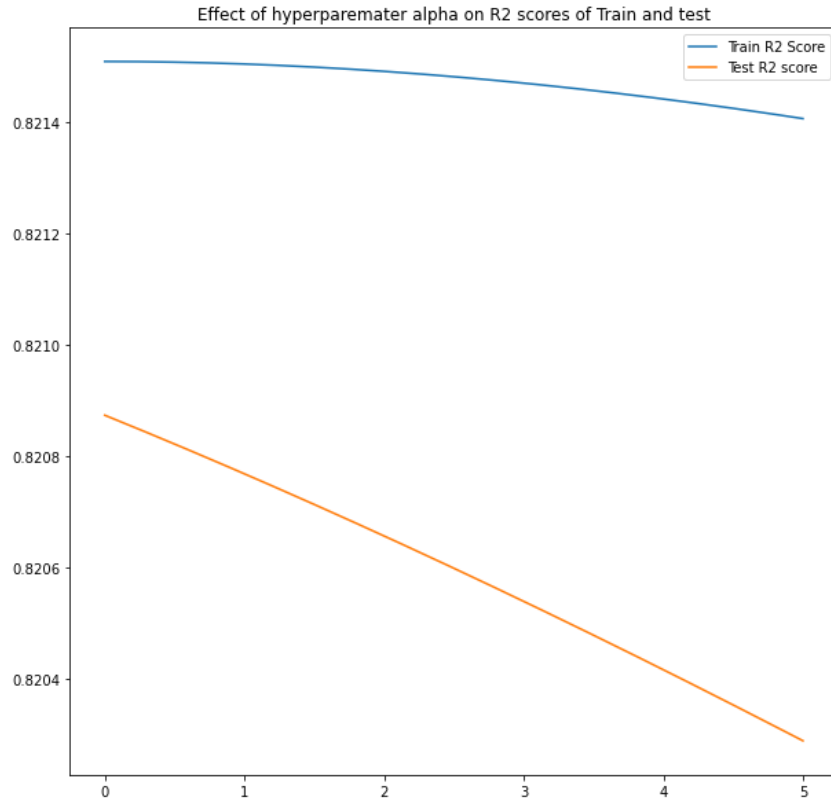
In [179]:

```
 1
 2  plt.figure(figsize = (10,10))
 3  sns.lineplot(lambdas,train_R2_score,)
 4  sns.lineplot(lambdas, test_R2_score)
 5  plt.legend(['Train R2 Score','Test R2 score'])
 6  plt.title("Effect of hyperparemater alpha on R2 scores of Train and test")
 7
 8
 9  plt.show()
10
```



In [180]:

```
1  RidgeModel = Ridge(alpha = 0.1)
2  RidgeModel.fit(X_train,y_train)
3  trainR2 = RidgeModel.score(X_train,y_train)
4  testR2 = RidgeModel.score(X_test,y_test)
```

In [181]:

```
1
2  trainR2,testR2
```

Out[181]:

(0.8215098726041209, 0.8208639536156423)

In [182]:

```
1
2  RidgeModel.coef_
```

Out[182]:

```
array([[0.02069489, 0.01929637, 0.00700953, 0.00298992, 0.01334235,
        0.07044884, 0.00987467]])
```

In [183]:

```
1  RidgeModel_coefs = pd.DataFrame(RidgeModel.coef_.reshape(1,-1),columns=df.columns[:-1])
2  RidgeModel_coefs["Intercept"] = RidgeModel.intercept_
3  RidgeModel_coefs
```

Out[183]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020695  | 0.019296    | 0.00701           | 0.00299 | 0.013342 | 0.070449 | 0.009875 | 0.722882 |

In [184]:

```
1  LinearRegression_Model_coefs
```

Out[184]:

| | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|---|---|---|---|---|---|---|---|
| **0** | 0.020675 | 0.019284 | 0.007001 | 0.002975 | 0.013338 | 0.070514 | 0.009873 | 0.722881 |

In [185]:

```python
1
2  y_pred = RidgeModel.predict(X_test)
3
4  print("MSE:",mean_squared_error(y_test,y_pred)) # MSE
5  print("RMSE:",np.sqrt(mean_squared_error(y_test,y_pred))) #RMSE
6  print("MAE :",mean_absolute_error(y_test,y_pred) ) # MAE
7  print("r2_score:",r2_score(y_test,y_pred)) # r2score
8  print("Adjusted R2 score :", AdjustedR2score(r2_score(y_test,y_pred),len(X),X.shape[1]))  # adjusted R2 score
9
```

```
MSE: 0.003459296191728331
RMSE: 0.058815781825359854
MAE : 0.040203055117056935
r2_score: 0.8208639536156423
Adjusted R2 score : 0.818315270028873
```

In [ ]:

```
1
```

In [186]:

```python
1
2  y_predicted = RidgeModel.predict(X_train)
3
4  residuals = (y_train - y_predicted)
5  plt.figure(figsize=(12,4))
6  plt.subplot(1,2,1)
7  sns.distplot(residuals)
8  plt.subplot(1,2,2)
9  stats.probplot(residuals.reshape(-1,), plot = plt)
10 plt.title('QQ plot for residual')
11 plt.show()
12
```



In [ ]:

```
1
```

# L1 regularization :

## Lasso :

In [189]:

```python
## Hyperparameter Tuning : for appropriate lambda value :

train_R2_score = []
test_R2_score = []
lambdas = []
train_test_difference_Of_R2 =  []
lambda_ = 0
while lambda_ <= 5:
    lambdas.append(lambda_)
    LassoModel = Lasso(alpha=lambda_)
    LassoModel.fit(X_train , y_train)
    trainR2 = LassoModel.score(X_train,y_train)
    testR2 = LassoModel.score(X_test,y_test)
    train_R2_score.append(trainR2)
    test_R2_score.append(testR2)

    lambda_ += 0.001
```

In [ ]:

```python

```

In [190]:

```python

plt.figure(figsize = (10,10))
sns.lineplot(lambdas,train_R2_score,)
sns.lineplot(lambdas, test_R2_score)
plt.legend(['Train R2 Score','Test R2 score'])
plt.title("Effect of hyperparemater alpha on R2 scores of Train and test")


plt.show()

```



In [ ]:

```python

```

In [191]:

```
1  LassoModel = Lasso(alpha=0.001)
2  LassoModel.fit(X_train , y_train)
3  trainR2 = LassoModel.score(X_train,y_train)
4  testR2 = LassoModel.score(X_test,y_test)
```

In [192]:

```
1  trainR2,testR2
```

Out[192]:

(0.82142983289567, 0.8198472607571161)

In [193]:

```
1
2  Lasso_Model_coefs = pd.DataFrame(LassoModel.coef_.reshape(1,-1),columns=df.columns[:-1])
3  Lasso_Model_coefs["Intercept"] = LassoModel.intercept_
4  Lasso_Model_coefs
```

Out[193]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020616  | 0.019069    | 0.006782          | 0.002808 | 0.012903 | 0.070605 | 0.009278 | 0.722863 |

In [194]:

```
1  RidgeModel_coefs
```

Out[194]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020695  | 0.019296    | 0.00701           | 0.00299 | 0.013342 | 0.070449 | 0.009875 | 0.722882 |

In [195]:

```
1  LinearRegression_Model_coefs
```

Out[195]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020675  | 0.019284    | 0.007001          | 0.002975 | 0.013338 | 0.070514 | 0.009873 | 0.722881 |

In [ ]:

```
1
```

In [196]:

```
1
2  y_predicted = LassoModel.predict(X_train)
3
4  residuals = (y_train - y_predicted)
5  plt.figure(figsize=(12,4))
6  plt.subplot(1,2,1)
7  sns.distplot(residuals)
8  plt.subplot(1,2,2)
9  stats.probplot(residuals.reshape(-1,), plot = plt)
10 plt.title('QQ plot for residual')
11 plt.show()
12
```
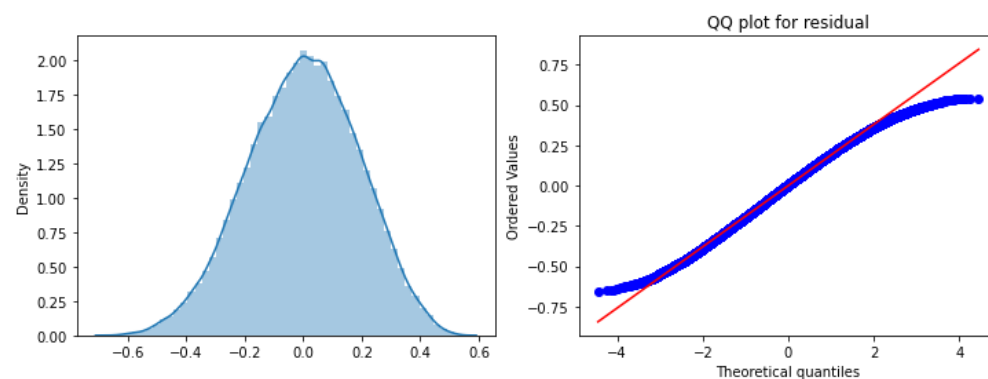
In [197]:

```
y_pred = LassoModel.predict(X_test)

print("MSE:",mean_squared_error(y_test,y_pred)) # MSE
print("RMSE:",np.sqrt(mean_squared_error(y_test,y_pred))) #RMSE
print("MAE :",mean_absolute_error(y_test,y_pred) ) # MAE
print("r2_score:",r2_score(y_test,y_pred)) # r2score
print("Adjusted R2 score :", AdjustedR2score(r2_score(y_test,y_pred),len(X),X.shape[1]))  # adjusted R2 score
```

```
MSE: 0.0034789295475193297
RMSE: 0.058982451182697807
MAE : 0.04022896061335951
r2_score: 0.8198472607571161
Adjusted R2 score : 0.8172841120280507
```

In [ ]:

```
```

# ElasticNet

## L1 and L2 regularisation :

**Elastic net linear regression uses the penalties from both the lasso and ridge techniques to regularize regression models.**

In [201]:

```
## Hyperparameter Tuning : for appropriate lambda value :

train_R2_score = []
test_R2_score = []
lambdas = []
train_test_difference_Of_R2 =  []
lambda_ = 0
while lambda_ <= 5:
    lambdas.append(lambda_)
    ElasticNet_model = ElasticNet(alpha=lambda_)
    ElasticNet_model.fit(X_train , y_train)
    trainR2 = ElasticNet_model.score(X_train,y_train)
    testR2 = ElasticNet_model.score(X_test,y_test)
    train_R2_score.append(trainR2)
    test_R2_score.append(testR2)

    lambda_ += 0.001
```

In [202]:

```python
plt.figure(figsize = (10,10))
sns.lineplot(lambdas,train_R2_score,)
sns.lineplot(lambdas, test_R2_score)
plt.legend(['Train R2 Score','Test R2 score'])
plt.title("Effect of hyperparemater alpha on R2 scores of Train and test")


plt.show()
```



In [203]:

```python
ElasticNet_model = ElasticNet(alpha=0.001)
ElasticNet_model.fit(X_train , y_train)
trainR2 = ElasticNet_model.score(X_train,y_train)
testR2 = ElasticNet_model.score(X_test,y_test)
```

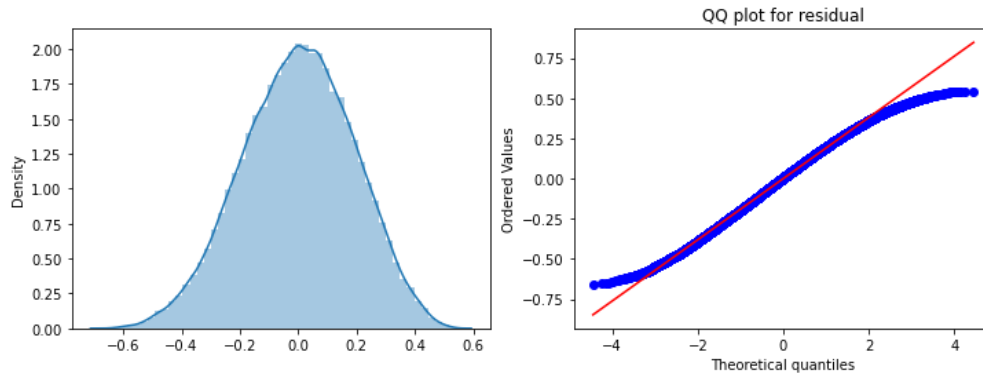In [204]:

```python

trainR2,testR2
```

Out[204]:

```
(0.8214893364453533, 0.8203602261096284)
```

In [205]:

```
1  y_predicted = ElasticNet_model.predict(X_train)
2
3  residuals = (y_train - y_predicted)
4  plt.figure(figsize=(12,4))
5  plt.subplot(1,2,1)
6  sns.distplot(residuals)
7  plt.subplot(1,2,2)
8  stats.probplot(residuals.reshape(-1,), plot = plt)
9  plt.title('QQ plot for residual')
10 plt.show()
11
```



In [206]:

```
1
2  y_pred = ElasticNet_model.predict(X_test)
3
4  print("MSE:",mean_squared_error(y_test,y_pred)) # MSE
5  print("RMSE:",np.sqrt(mean_squared_error(y_test,y_pred))) #RMSE
6  print("MAE :",mean_absolute_error(y_test,y_pred) ) # MAE
7  print("r2_score:",r2_score(y_test,y_pred)) # r2score
8  print("Adjusted R2 score :", AdjustedR2score(r2_score(y_test,y_pred),len(X),X.shape[1]))  # adjusted R2 score
9
```

```
MSE: 0.003469023673596966
RMSE: 0.058898418260569324
MAE : 0.04021407699792928
r2_score: 0.8203602261096284
Adjusted R2 score : 0.8178043756680987
```

In [207]:

```
1
2  ElasticNet_model_coefs = pd.DataFrame(ElasticNet_model.coef_.reshape(1,-1),columns=df.columns[:-1])
3  ElasticNet_model_coefs["Intercept"] = ElasticNet_model.intercept_
4  ElasticNet_model_coefs
5
```

Out[207]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020679 | 0.019199 | 0.006908 | 0.00292 | 0.013128 | 0.070437 | 0.009581 | 0.722873 |

In [208]:

```
1  RidgeModel_coefs
```

Out[208]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020695 | 0.019296 | 0.00701 | 0.00299 | 0.013342 | 0.070449 | 0.009875 | 0.722882 |

In [209]:

```
1  Lasso_Model_coefs
```

Out[209]:

|   | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|-----------|-------------|-------------------|-----|-----|------|----------|-----------|
| 0 | 0.020616 | 0.019069 | 0.006782 | 0.002808 | 0.012903 | 0.070605 | 0.009278 | 0.722863 |

In [210]:

```
1
2  LinearRegression_Model_coefs
```

Out[210]:

| | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept |
|---|---|---|---|---|---|---|---|---|
| **0** | 0.020675 | 0.019284 | 0.007001 | 0.002975 | 0.013338 | 0.070514 | 0.009873 | 0.722881 |

In [ ]:

```
1
```

In [211]:

```
1
2  y_pred = ElasticNet_model.predict(X_test)
3  ElasticNet_model_metrics = []
4  ElasticNet_model_metrics.append(mean_squared_error(y_test,y_pred)) # MSE
5  ElasticNet_model_metrics.append(np.sqrt(mean_squared_error(y_test,y_pred))) #RMSE
6  ElasticNet_model_metrics.append(mean_absolute_error(y_test,y_pred) ) # MAE
7  ElasticNet_model_metrics.append(r2_score(y_test,y_pred)) # r2score
8  ElasticNet_model_metrics.append(AdjustedR2score(r2_score(y_test,y_pred),len(X),X.shape[1]))  # adjusted R2 score
9
```

In [ ]:

```
1
```

In [212]:

```
1  y_pred = LinearRegression.predict(X_test)
2  LinearRegression_model_metrics = []
3  LinearRegression_model_metrics.append(mean_squared_error(y_test,y_pred)) # MSE
4  LinearRegression_model_metrics.append(np.sqrt(mean_squared_error(y_test,y_pred))) #RMSE
5  LinearRegression_model_metrics.append(mean_absolute_error(y_test,y_pred) ) # MAE
6  LinearRegression_model_metrics.append(r2_score(y_test,y_pred)) # r2score
7  LinearRegression_model_metrics.append(AdjustedR2score(r2_score(y_test,y_pred),len(X),X.shape[1]))  # adjusted R2 score
8
```

In [ ]:

```
1
```

In [213]:

```
1
2  y_pred = RidgeModel.predict(X_test)
3  RidgeModel_model_metrics = []
4  RidgeModel_model_metrics.append(mean_squared_error(y_test,y_pred)) # MSE
5  RidgeModel_model_metrics.append(np.sqrt(mean_squared_error(y_test,y_pred))) #RMSE
6  RidgeModel_model_metrics.append(mean_absolute_error(y_test,y_pred) ) # MAE
7  RidgeModel_model_metrics.append(r2_score(y_test,y_pred)) # r2score
8  RidgeModel_model_metrics.append(AdjustedR2score(r2_score(y_test,y_pred),len(X),X.shape[1]))  # adjusted R2 score
9
```

In [ ]:

```
1
```

In [214]:

```
1
2  y_pred = LassoModel.predict(X_test)
3  LassoModel_model_metrics = []
4  LassoModel_model_metrics.append(mean_squared_error(y_test,y_pred)) # MSE
5  LassoModel_model_metrics.append(np.sqrt(mean_squared_error(y_test,y_pred))) #RMSE
6  LassoModel_model_metrics.append(mean_absolute_error(y_test,y_pred) ) # MAE
7  LassoModel_model_metrics.append(r2_score(y_test,y_pred)) # r2score
8  LassoModel_model_metrics.append(AdjustedR2score(r2_score(y_test,y_pred),len(X),X.shape[1]))  # adjusted R2 score
9
```

In [ ]:

```
1
```

In [215]:

```
1  ElasticNet_model_metrics
```

Out[215]:

```
[0.003469023673596966,
 0.058898418260569324,
 0.04021407699792928,
 0.8203602261096284,
 0.8178043756680987]
```

In [ ]:

```
1
```

In [216]:

```
1
2  A = pd.DataFrame([LinearRegression_model_metrics,LassoModel_model_metrics,RidgeModel_model_metrics,ElasticNet_model_metrics],colu
3  A
```

Out[216]:

|  | MSE | RMSE | MAE | R2_SCORE | ADJUSTED_R2 |
|---|---|---|---|---|---|
| Linear Regression Model | 0.003459 | 0.058814 | 0.040200 | 0.820874 | 0.818326 |
| Lasso Regression Model | 0.003479 | 0.058982 | 0.040229 | 0.819847 | 0.817284 |
| Ridge Regression Model | 0.003459 | 0.058816 | 0.040203 | 0.820864 | 0.818315 |
| ElasticNet Regression Model | 0.003469 | 0.058898 | 0.040214 | 0.820360 | 0.817804 |

In [217]:

```
1  B = pd.DataFrame(LinearRegression_Model_coefs.append(Lasso_Model_coefs).append(RidgeModel_coefs).append(ElasticNet_model_coefs))
2  B.index = ["Linear Regression Model","Lasso Regression Model","Ridge Regression Model","ElasticNet Regression Model"]
3
```

In [218]:

```
1  REPORT = B.reset_index().merge(A.reset_index())
```

In [219]:

```
1
2  REPORT = REPORT.set_index("index")
3  REPORT
```
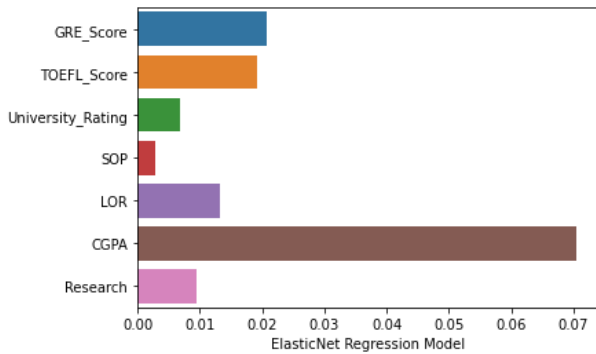
Out[219]:

|  | GRE_Score | TOEFL_Score | University_Rating | SOP | LOR | CGPA | Research | Intercept | MSE | RMSE | MAE | R2_SCOR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **index** | | | | | | | | | | | | |
| Linear Regression Model | 0.020675 | 0.019284 | 0.007001 | 0.002975 | 0.013338 | 0.070514 | 0.009873 | 0.722881 | 0.003459 | 0.058814 | 0.040200 | 0.82087 |
| Lasso Regression Model | 0.020616 | 0.019069 | 0.006782 | 0.002808 | 0.012903 | 0.070605 | 0.009278 | 0.722863 | 0.003479 | 0.058982 | 0.040229 | 0.81984 |
| Ridge Regression Model | 0.020695 | 0.019296 | 0.007010 | 0.002990 | 0.013342 | 0.070449 | 0.009875 | 0.722882 | 0.003459 | 0.058816 | 0.040203 | 0.82086 |
| ElasticNet Regression Model | 0.020679 | 0.019199 | 0.006908 | 0.002920 | 0.013128 | 0.070437 | 0.009581 | 0.722873 | 0.003469 | 0.058898 | 0.040214 | 0.82036 |

In [221]:

```
1  sns.barplot(y  = REPORT.loc["ElasticNet Regression Model"][0:7].index,
2             x = REPORT.loc["ElasticNet Regression Model"][0:7])
```

Out[221]:

```
<AxesSubplot:xlabel='ElasticNet Regression Model'>
```



In [ ]:

```
1
```

In [ ]:

```
1
```

# Insights , Feature Importance and Interpretations and Recommendations :

fist column was observed as unique row identifier which was dropped and was not required for model building.

University Rating , SOP and LOR strength and research are seems to be discrete random Variables , but also ordinal numeric data.

all the other features are numeric, ordinal and continuous.

No null values were present in data.

No Significant amount of outliers were found in data.

Chance of admission(target variable) and GRE score(an independent feature) are nearly normally distrubted.

Independent Variables (Input data): GRE Score, TOEFL Score, University Rating, SOP, LOR, CGPA, Research

Target/Dependent Variable : Chance of Admit (the value we want to predict)

from correlation heatmap , we can observe GRE score, TOEFL score and CGPA have very high correlation with Change of admission.

University rating, SOP ,LOR and Research have comparatively slightly less correlated than other features.

chances of admit is a probability measure , which is within 0 to 1 which is good (no outliers or missleading data in column).

Range of GRE score looks like between 290 to 340.

range of TOEFL score is between 92 to 120.

university rating , SOP and LOR are distributed between range of 1 to 5.

CGPA range is between 6.8 to 9.92.

From boxplots (distribution of chance of admition (probability of getting admition) as per GRE score ) : with higher GRE score , there is high probability of getting an admition .

Students having high toefl score , has higher probability of getting admition .

from count plots, we can observe , statement of purpose SOP strength is positively correlated with Chance of Admission .

we can also similar pattern in Letter of Recommendation Stength and University rating , have positive correlation with Chaces of Admission .

Student having research has higher chances of Admission , but also we can observe some outliers within that caregory.

In [ ]:

```
1
```

# Actionable Insights and Recommendations :

education institute can not just help student to improve their CGPA score but also assist them writing good LOR and SOP thus helping them admit to better university.

The education institute can not just help student to improve their GRE Score but can also assist them writing good LOR and SOP thus helping them admit to a better University.

Awareness of CGPA and Reserach Capabilities : Seminars can be organised to increase the awareness regarding CGPA and Research Capablities to enhance the chance of admit.

Any student can never change their current state of attributes so awareness and marketing campaign need to surveyed hence creating a first impression on student at undergraduate level, which wont just increase company's popularity but will also help sudent get prepared for future plans in advance.

A dashboard can be created for students whenever they loged in into your website, hence allowing a healthy competition also to create a progress report for students.

Additional features like number of hours they put in studing, watching lectures, assignments soved percentage, marks in mock test can result a better report for every student to judge themselves and improve on their own.

# Regression Analysis :

from regression analysis (above bar chart and REPORT file), we can observe the CGPA is the most Important feature for prediciing the chances of admission.

other important features are GRE and TOEFL score .

after first Regression Model, checked for Multicolinearity . Getting all the VIF scores below 5 , showing there's no high multicolinearity.

all the residuals are not perfectly normally distributed. and so residual plot we can observe some level of heteroscedasticity.

regularised model ridge and lasso both give very similar results to Linear Regression Model.

similarly ElasticNet (L1+L2) also returns very similar results. along with rest of all the model metrics.