

# Importing basic libraries

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 %matplotlib inline
        5 import seaborn as sns
        6 import warnings
        7 warnings.filterwarnings('ignore')
        8 from sklearn.model_selection import train_test_split
        9 from sklearn.preprocessing import StandardScaler
       10 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
       11 from sklearn.svm import SVR
       12 from sklearn.linear_model import LinearRegression, Ridge, Lasso, ElasticNet
```

<frozen importlib.\_bootstrap>:228: RuntimeWarning: scipy.\_lib.messagestream.MessageStream size changed, may indicate binary incompatibility. Expected 56 from C header, got 64 from PyObject

```
In [2]: 1 # Reading the data
        2 df = pd.read_csv('https://raw.githubusercontent.com/Jhoie/Admission-Prediction/main/Admission_Predict.csv')
```

```
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]: (400, 9)

```
In [5]: 1 df.nunique()
```

Out[5]: Serial No. 400  
GRE Score 49  
TOEFL Score 29  
University Rating 5  
SOP 9  
LOR 9  
CGPA 168  
Research 2  
Chance of Admit 60  
dtype: int64

```
In [6]: 1 df.columns
```

Out[6]: Index(['Serial No.', 'GRE Score', 'TOEFL Score', 'University Rating', 'SOP',  
'LOR ', 'CGPA', 'Research', 'Chance of Admit '],  
dtype='object')

```
In [7]: 1 # Removing the extra space in the column name
        2
        3 df.columns = [x.strip() for x in df.columns]
        4 df.columns
```

Out[7]: Index(['Serial No.', 'GRE Score', 'TOEFL Score', 'University Rating', 'SOP',  
'LOR', 'CGPA', 'Research', 'Chance of Admit'],  
dtype='object')

```
In [8]: 1 # Checking for null values
        2 df.isnull().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

# Statistical Analysis

In [9]:

1df.describe().T

Out[9]:

	count	mean	std	min	25%	50%	75%	max
Serial No.	400.0	200.500000	115.614301	1.00	100.75	200.50	300.2500	400.00
GRE Score	400.0	316.807500	11.473646	290.00	308.00	317.00	325.0000	340.00
TOEFL Score	400.0	107.410000	6.069514	92.00	103.00	107.00	112.0000	120.00
University Rating	400.0	3.087500	1.143728	1.00	2.00	3.00	4.0000	5.00
SOP	400.0	3.400000	1.006869	1.00	2.50	3.50	4.0000	5.00
LOR	400.0	3.452500	0.898478	1.00	3.00	3.50	4.0000	5.00
CGPA	400.0	8.598925	0.596317	6.80	8.17	8.61	9.0625	9.92
Research	400.0	0.547500	0.498362	0.00	0.00	1.00	1.0000	1.00
Chance of Admit	400.0	0.724350	0.142609	0.34	0.64	0.73	0.8300	0.97

In [10]:

1df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Serial No.            400 non-null   int64
1   GRE Score             400 non-null   int64
2   TOEFL Score           400 non-null   int64
3   University Rating     400 non-null   int64
4   SOP                   400 non-null   float64
5   LOR                   400 non-null   float64
6   CGPA                  400 non-null   float64
7   Research              400 non-null   int64
8   Chance of Admit       400 non-null   float64
dtypes: float64(4), int64(5)
memory usage: 28.2 KB
```

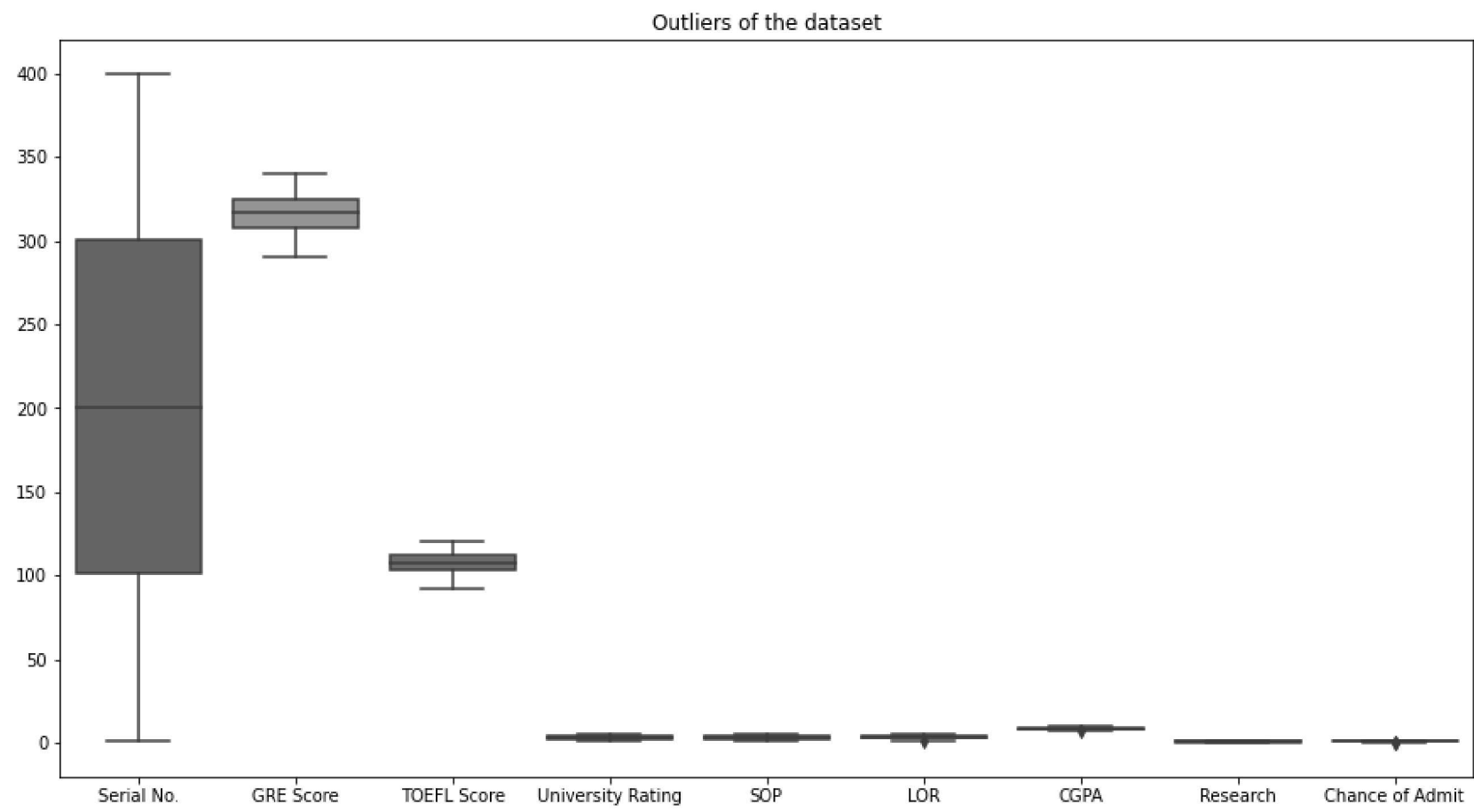
In [11]:

1## Univariate Analysis
2for i in df.columns:
3 print('-----')
4 print(i)
5 print('-----')
6 print(df[i].value\_counts())

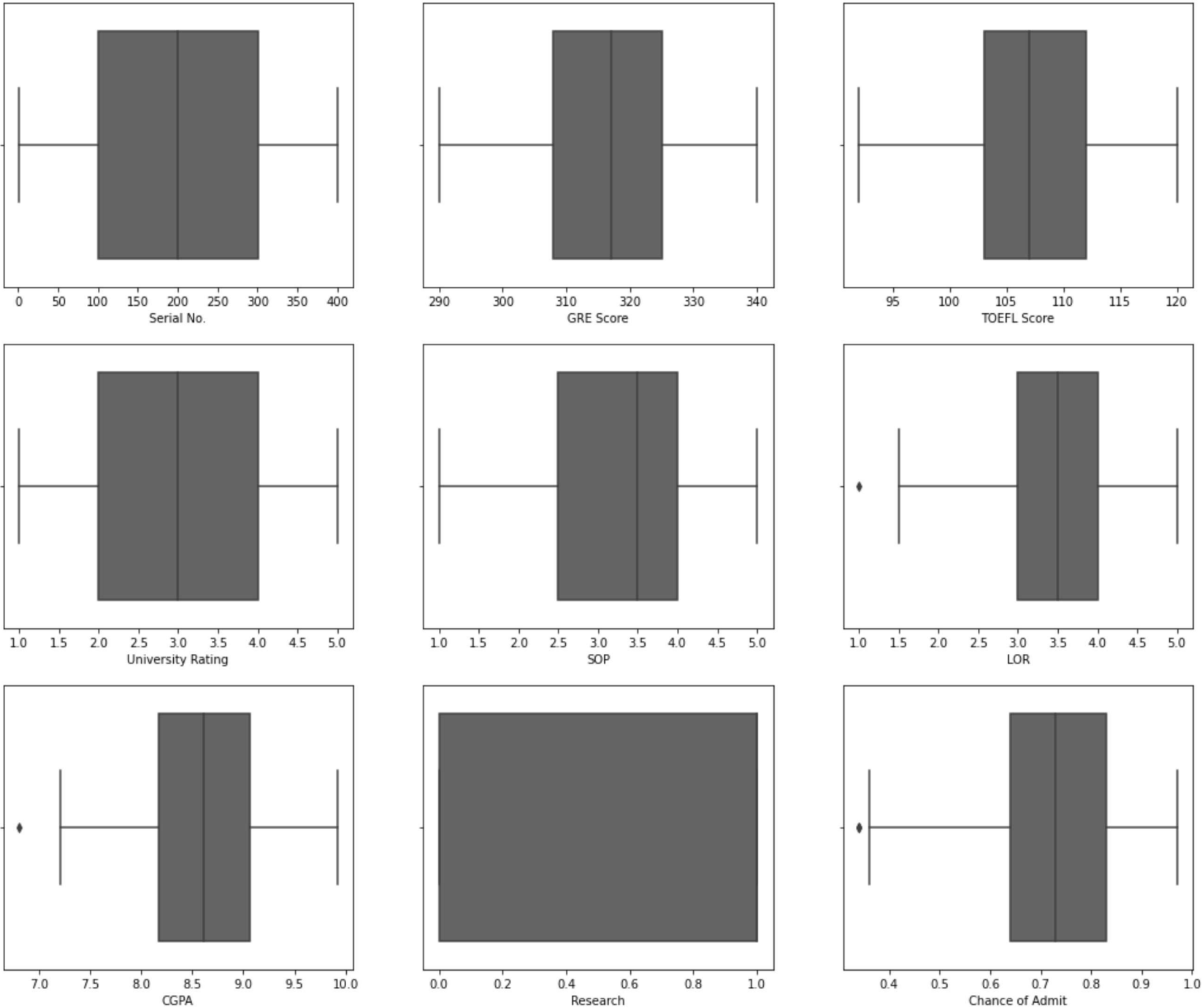
```
336    5
296    5
303    5
302    5
335    4
295    4
338    4
297    4
333    4
339    3
294    2
290    2
337    1
293    1
Name: GRE Score, dtype: int64
-----
TOEFL Score
-----
110    37
105    28
```

```
In [12]: 1 # checking the outlier
        2
        3 plt.figure(figsize=(15,8))
        4 sns.boxplot(data = df)
        5 plt.title('Outliers of the dataset')
```

Out[12]: Text(0.5, 1.0, 'Outliers of the dataset')

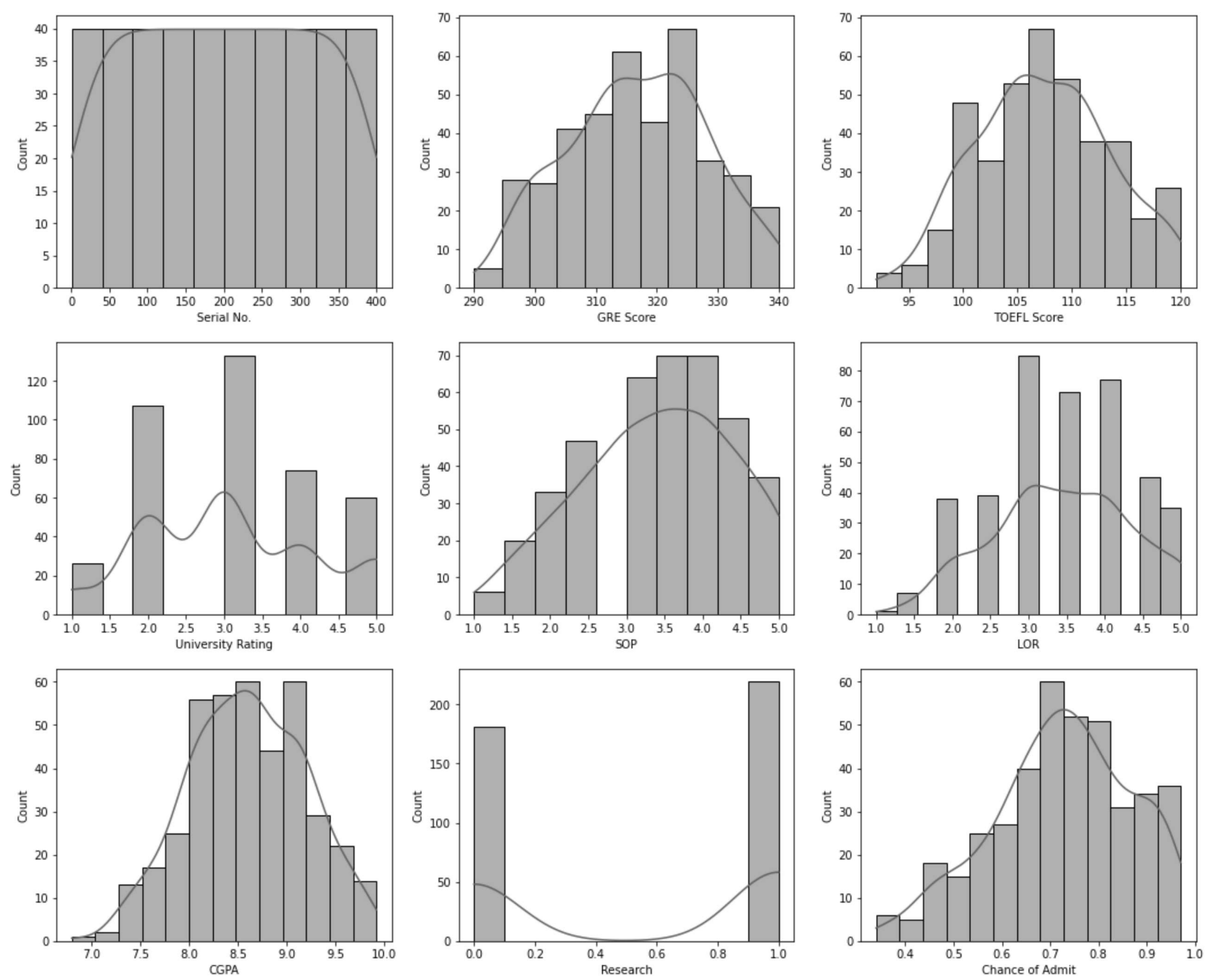


```
In [13]: 1 plt.figure(figsize = (18,15))
2 plotnumber=1
3 for i in df.columns:
4     if plotnumber<=9:
5         ax = plt.subplot(3,3,plotnumber)
6         sns.boxplot(df[i])
7         plotnumber+=1
8 plt.show()
9
10 '''
11
12 plt.figure(figsize=(18,15))
13 plotnumber=1
14 for i in df.columns:
15     if plotnumber<=6:
16         ax = plt.subplot(3,2,plotnumber)
17         sns.boxplot(df[i])
18         plotnumber+=1
19 plt.show()
20 '''
```

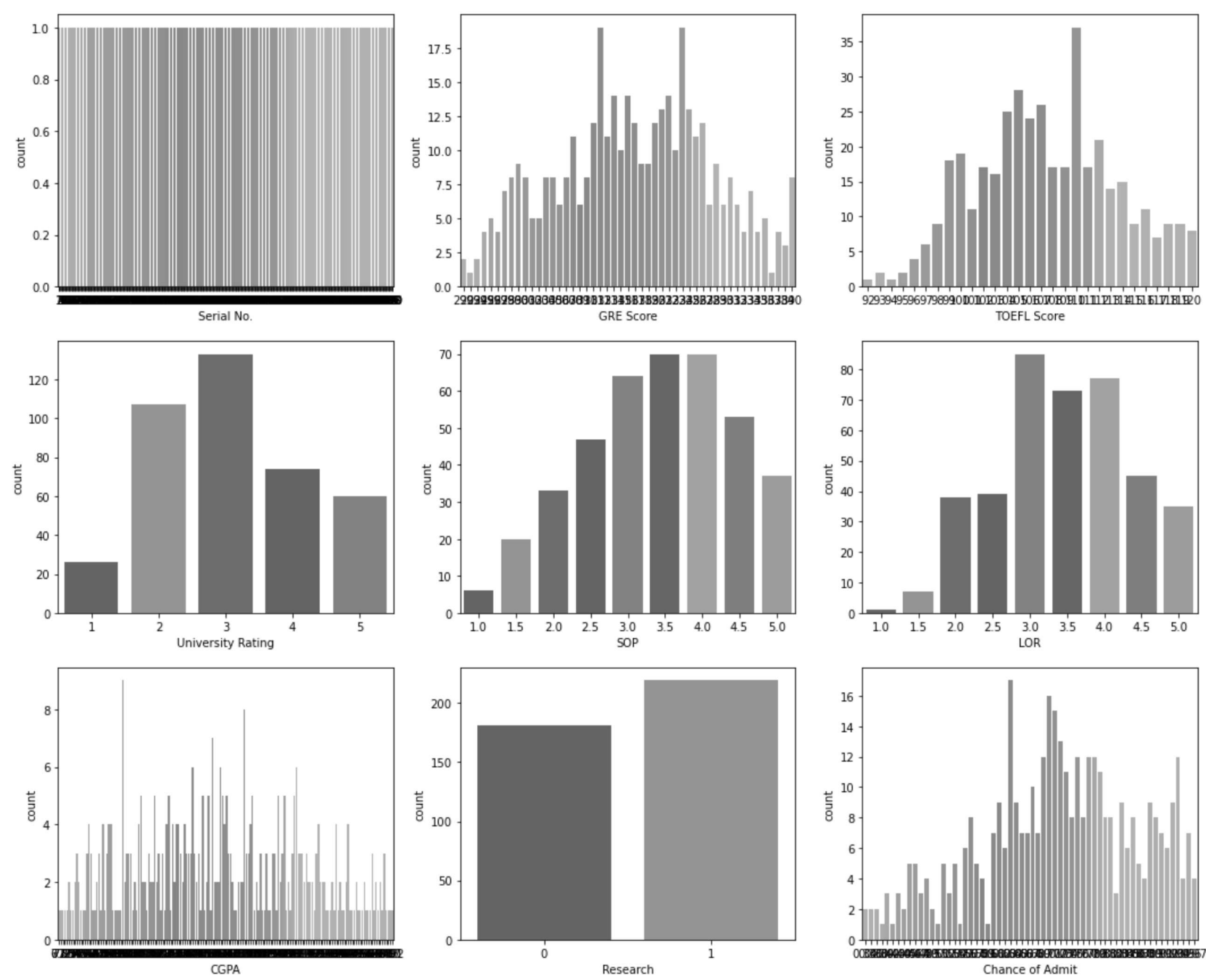


```
Out[13]: '\n\nplt.figure(figsize=(18,15))\n\nplotnumber=1\n\nfor i in df.columns:\n\n    if plotnumber<=6:\n\n        ax = pl\n\n        t.subplot(3,2,plotnumber)\n\n        sns.boxplot(df[i])\n\n        plotnumber+=1\n\nplt.show()\n'
```

```
In [14]: 1 # Getting histplot of the features
2
3 plt.figure(figsize=(18,15))
4 plotnumber=1
5 for i in df.columns:
6     if plotnumber<=9:
7         ax=plt.subplot(3,3,plotnumber)
8         sns.histplot(df[i],kde=True)
9         plotnumber+=1
10 plt.show()
```

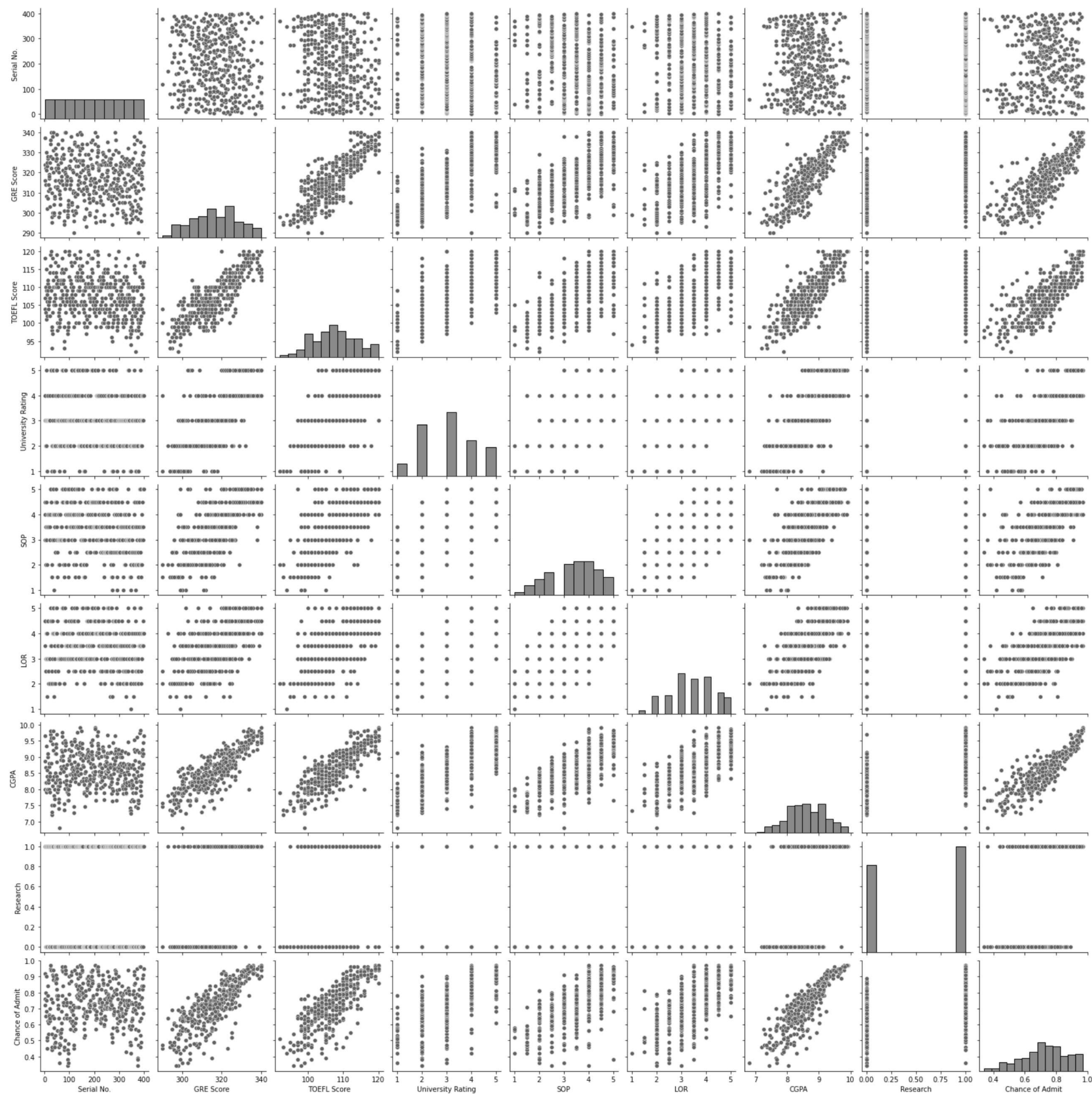


```
In [15]: 1 # Getting countplot of the features
2
3 plt.figure(figsize=(18,15))
4 plotnumber=1
5 for i in df.columns:
6     if plotnumber<=9:
7         ax=plt.subplot(3,3,plotnumber)
8         sns.countplot(df[i])
9         plotnumber+=1
10 plt.show()
```



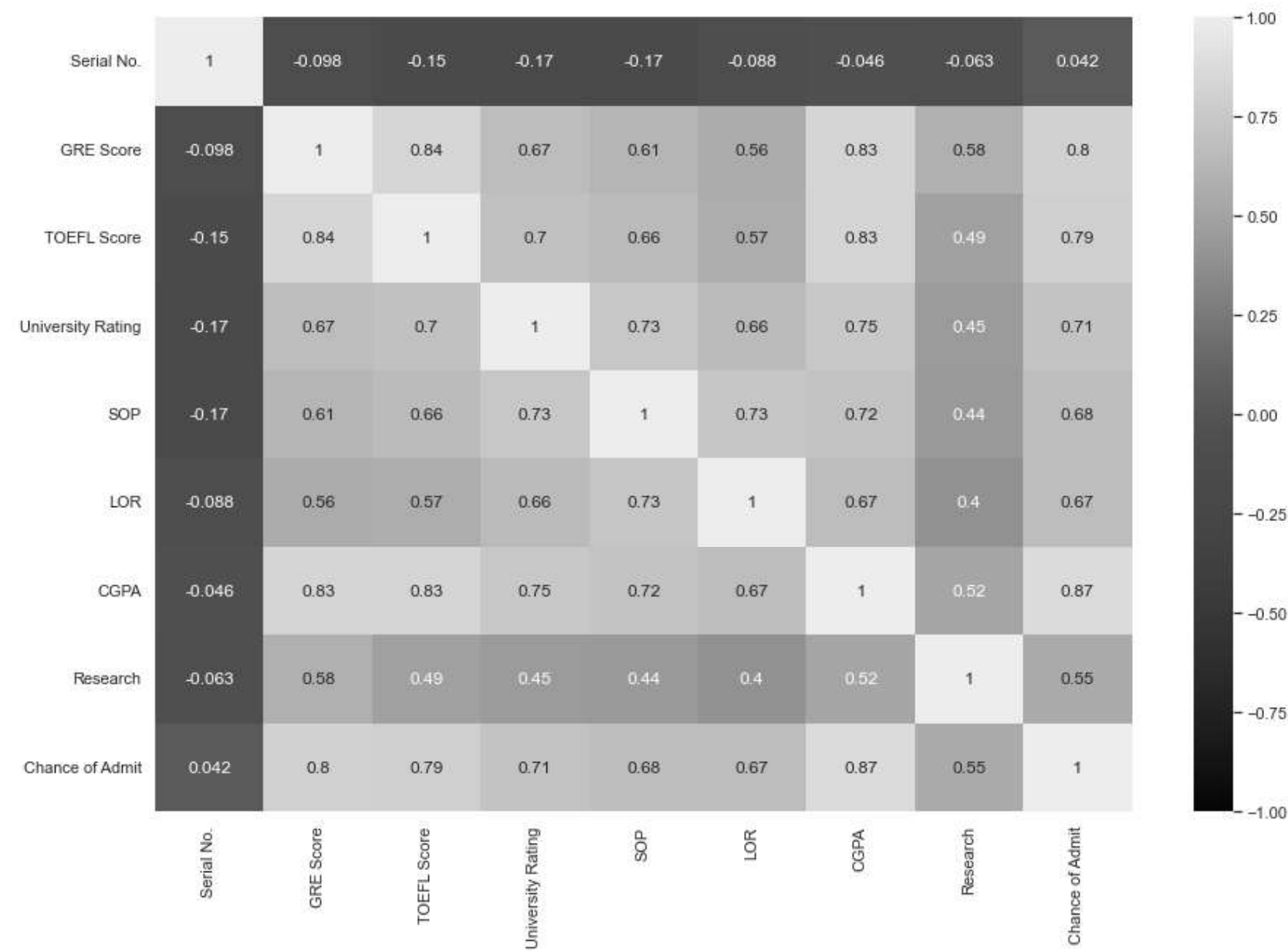
```
In [16]: 1 ## Bivariate Analysis
        2 sns.pairplot(df)
```

Out[16]: <seaborn.axisgrid.PairGrid at 0x1ef5e7d5850>



```
In [17]: 1 sns.set(rc={'figure.figsize':(15,10)})
2 sns.heatmap(data=df.corr(),annot=True, vmin=-1, vmax=1)
```

Out[17]: <AxesSubplot:>



```
1 Observation:- Drop 'Serial No.' feature from the dataset, because it is of no use.
2 GRE Score and TOFEL Score are highly correlated with CGPA
```

```
In [18]: 1 df.columns
```

Out[18]: Index(['Serial No.', 'GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR', 'CGPA', 'Research', 'Chance of Admit'], dtype='object')

```
In [19]: 1 # df.drop(['Serial No.','GRE Score','TOEFL Score'],axis=1,inplace=True)
```

```
In [20]: 1 ## Segregating dependent and independent feature
2
3 X = df.iloc[ : , :-1]
4 y = df.iloc[ : , -1]
```

```
In [21]: 1 X.head()
```

Out[21]:

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

## Train Test Split

```
In [22]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.33,random_state=4)
```

```
1 # Checking VIF
2
3 from statsmodels.stats.outliers_influence import variance_inflation_factor
```



```
In [23]: 1 # Scaling
2 sc = StandardScaler()
3
4 X_train_tf = sc.fit_transform(X_train)
5 X_test_tf = sc.transform(X_test) # on;y transform to avoid data Leakage
```

## Model Creation

```
In [24]: 1 # SVR Model
2 model_svr = SVR()
```

```
In [25]: 1 model_svr.fit(X_train_tf,y_train)
```

Out[25]:

▼ SVR

SVR()

```
In [26]: 1 ## Train Accuracy
2 train_score = model_svr.score(X_train_tf,y_train)
3 train_score
```

Out[26]: 0.7947034504013842

```
In [27]: 1 # Store model using pickle
2 import pickle
3
4 with open('AdmissionSVR.pkl','wb') as f:
5     pickle.dump(model_svr,f)
```

```
In [28]: 1 model_svr_load = pickle.load(open('AdmissionSVR.pkl','rb'))
```

```
In [29]: 1 y_pred_svr = model_svr.predict(X_test_tf)
```

```
In [30]: 1 print('MSE : ',round(mean_squared_error(y_test,y_pred_svr),2))
2 print('MAE : ',round(mean_absolute_error(y_test,y_pred_svr),2))
```

MSE : 0.0  
MAE : 0.06

```
In [31]: 1 svr_r2_score = r2_score(y_test, y_pred_svr)
2 svr_adj_r2_score = 1 - ((1-svr_r2_score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1))
```

```
In [32]: 1 print('R-squared Accuracy : ',round(svr_r2_score*100,2))
2 print('Adjusted R-squared Accuracy : ',round(svr_adj_r2_score*100,2))
```

R-squared Accuracy : 76.61  
Adjusted R-squared Accuracy : 75.08

```
In [33]: 1 ## GridSearchCV
2
3 from sklearn.model_selection import GridSearchCV
```

```
In [34]: 1 param_grid = {'C': [0.1,1, 10, 100], 'gamma': [1,0.1,0.01,0.001], 'kernel':['linear','rbf']}
```

```
In [35]: 1 model_grid_svr = GridSearchCV(SVR(), param_grid, verbose=3)
2 model_grid_svr.fit(X_train_tf, y_train)
[CV 2/5] END .....C=0.1, gamma=0.01, kernel=rbf;, score=0.793 total time= 0.0s
[CV 3/5] END .....C=0.1, gamma=0.01, kernel=rbf;, score=0.718 total time= 0.0s
[CV 4/5] END .....C=0.1, gamma=0.01, kernel=rbf;, score=0.659 total time= 0.0s
[CV 5/5] END .....C=0.1, gamma=0.01, kernel=rbf;, score=0.651 total time= 0.0s
[CV 1/5] END .C=0.1, gamma=0.001, kernel=linear;, score=0.848 total time= 0.0s
[CV 2/5] END .C=0.1, gamma=0.001, kernel=linear;, score=0.866 total time= 0.0s
[CV 3/5] END .C=0.1, gamma=0.001, kernel=linear;, score=0.718 total time= 0.0s
[CV 4/5] END .C=0.1, gamma=0.001, kernel=linear;, score=0.678 total time= 0.0s
[CV 5/5] END .C=0.1, gamma=0.001, kernel=linear;, score=0.650 total time= 0.0s
[CV 1/5] END ....C=0.1, gamma=0.001, kernel=rbf;, score=0.600 total time= 0.0s
[CV 2/5] END ....C=0.1, gamma=0.001, kernel=rbf;, score=0.572 total time= 0.0s
[CV 3/5] END ....C=0.1, gamma=0.001, kernel=rbf;, score=0.599 total time= 0.0s
[CV 4/5] END ....C=0.1, gamma=0.001, kernel=rbf;, score=0.475 total time= 0.0s
[CV 5/5] END ....C=0.1, gamma=0.001, kernel=rbf;, score=0.569 total time= 0.0s
[CV 1/5] END .....C=1, gamma=1, kernel=linear;, score=0.848 total time= 0.0s
[CV 2/5] END .....C=1, gamma=1, kernel=linear;, score=0.866 total time= 0.0s
[CV 3/5] END .....C=1, gamma=1, kernel=linear;, score=0.718 total time= 0.0s
[CV 4/5] END .....C=1, gamma=1, kernel=linear;, score=0.673 total time= 0.0s
[CV 5/5] END .....C=1, gamma=1, kernel=linear;, score=0.650 total time= 0.0s
[CV 1/5] END .....C=1, gamma=1, kernel=rbf;, score=0.251 total time= 0.0s
```

```
In [36]: 1 print(model_grid_svr.best_estimator_)
SVR(C=0.1, gamma=1, kernel='linear')
```

```
In [37]: 1 # Train Accuracy
2
3 model_grid_svr.score(X_train_tf,y_train)
```

Out[37]: 0.8027928629188807

```
In [38]: 1 # Test Accuracy
2
3 y_pred_grid = model_grid_svr.predict(X_test_tf)
```

```
In [39]: 1 print('MSE : ',round(mean_squared_error(y_test,y_pred_grid),2))
2 print('MAE : ',round(mean_absolute_error(y_test,y_pred_grid),2))
```

MSE : 0.0  
MAE : 0.05

```
In [40]: 1 grid_svr_r2_score = r2_score(y_test, y_pred_grid)
2 grid_svr_adj_r2_score = 1 - ((1-grid_svr_r2_score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1))
```

```
In [41]: 1 print('R-squared Accuracy : ',round(grid_svr_r2_score*100,2))
2 print('Adjusted R-squared Accuracy : ',round(grid_svr_adj_r2_score*100,2))
```

R-squared Accuracy : 80.07  
Adjusted R-squared Accuracy : 78.78

```
In [42]: 1 ## Linear Regression Model
2
3 lin_reg = LinearRegression()
4 lin_reg
```

Out[42]: ▾ LinearRegression  
LinearRegression()

```
In [43]: 1 lin_reg.fit(X_train_tf, y_train)
```

Out[43]: ▾ LinearRegression  
LinearRegression()

```
In [44]: 1 # Train Accuracy
2
3 lin_reg.score(X_train_tf, y_train)
```

Out[44]: 0.8168014873134735

```
In [45]: 1 lin_reg_pred = lin_reg.predict(X_test_tf)
```

```
In [46]: 1 print('MSE : ',round(mean_squared_error(y_test,lin_reg_pred),2))
2 print('MAE : ',round(mean_absolute_error(y_test,lin_reg_pred),2))
```

MSE : 0.0  
MAE : 0.05

```
In [47]: 1 lin_reg_r2_score = r2_score(y_test, lin_reg_pred)
2 lin_reg_adj_r2_score = 1 - ((1-lin_reg_r2_score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1))
```

```
In [48]: 1 print('R-squared Accuracy : ',round(lin_reg_r2_score*100,2))
2 print('Adjusted R-squared Accuracy : ',round(lin_reg_adj_r2_score*100,2))
```

R-squared Accuracy : 81.09  
Adjusted R-squared Accuracy : 79.86

```
In [49]: 1 # Ridge Regression Model
2
3 ridge_reg = Ridge()
4 ridge_reg
```

Out[49]: ▾ Ridge  
Ridge()

```
In [50]: 1 ridge_reg.fit(X_train_tf, y_train)

Out[50]: ▾ Ridge
          Ridge()

In [51]: 1 # Train Accuracy
          2 ridge_reg.score(X_train_tf, y_train)

Out[51]: 0.8167919100113752

In [52]: 1 ridge_reg_pred = ridge_reg.predict(X_test_tf)

In [53]: 1 print('MSE : ',round(mean_squared_error(y_test,ridge_reg_pred),2))
          2 print('MAE : ',round(mean_absolute_error(y_test,ridge_reg_pred),2))

MSE : 0.0
MAE : 0.05

In [54]: 1 ridge_reg_r2_score = r2_score(y_test, ridge_reg_pred)
          2 ridge_reg_adj_r2_score = 1 - ((1-ridge_reg_r2_score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1))

In [55]: 1 print('R-squared Accuracy : ',round(ridge_reg_r2_score*100,2))
          2 print('Adjusted R-squared Accuracy : ',round(ridge_reg_adj_r2_score*100,2))

R-squared Accuracy : 81.11
Adjusted R-squared Accuracy : 79.88

In [56]: 1 ## Lasso Regressionj Model
          2 lasso_reg = Lasso()
          3 lasso_reg.fit(X_train_tf, y_train)

Out[56]: ▾ Lasso
          Lasso()

In [57]: 1 lasso_reg_pred = lasso_reg.predict(X_test_tf)

In [58]: 1 print('MSE : ',round(mean_squared_error(y_test,lasso_reg_pred),2))
          2 print('MAE : ',round(mean_absolute_error(y_test,lasso_reg_pred),2))

MSE : 0.02
MAE : 0.12

In [59]: 1 lasso_reg_r2_score = r2_score(y_test, lasso_reg_pred)
          2 lasso_reg_adj_r2_score = 1 - ((1-lasso_reg_r2_score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1))

In [60]: 1 print('R-squared Accuracy : ',round(lasso_reg_r2_score*100,2))
          2 print('Adjusted R-squared Accuracy : ',round(lasso_reg_adj_r2_score*100,2))

R-squared Accuracy : -1.32
Adjusted R-squared Accuracy : -7.91

In [61]: 1 ## Elastic Net Regression Model
          2
          3 elastic_reg = ElasticNet()
          4 elastic_reg

Out[61]: ▾ ElasticNet
          ElasticNet()

In [62]: 1 elastic_reg.fit(X_train_tf, y_train)

Out[62]: ▾ ElasticNet
          ElasticNet()

In [63]: 1 # Train Accuracy
          2 elastic_reg.score(X_train_tf,y_train )

Out[63]: 0.0

In [64]: 1 elastic_reg_pred = elastic_reg.predict(X_test_tf)
```

```
In [65]: 1 print('MSE : ',round(mean_squared_error(y_test,elastic_reg_pred),2))
2 print('MAE : ',round(mean_absolute_error(y_test,elastic_reg_pred),2))

MSE :  0.02
MAE :  0.12

In [66]: 1 elastic_reg_r2_score = r2_score(y_test, elastic_reg_pred)
2 elastic_reg_adj_r2_score = 1 - ((1-elastic_reg_r2_score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1))

In [67]: 1 print('R-squared Accuracy : ',round(elastic_reg_r2_score*100,2))
2 print('Adjusted R-squared Accuracy : ',round(elastic_reg_adj_r2_score*100,2))

R-squared Accuracy :  -1.32
Adjusted R-squared Accuracy :  -7.91
```

## Comparison of all Models

```
In [68]: 1 print('-----')
2 print('MSE : \n1. SVR : ',round(mean_squared_error(y_test, y_pred_svr),2))
3 print('2. Linear Regression : ',round(mean_squared_error(y_test,lin_reg_pred),2))
4 print('3. Ridge Regression : ',round(mean_squared_error(y_test, ridge_reg_pred),2))
5 print('4. Elastic Net Regression : ',round(mean_squared_error(y_test, elastic_reg_pred),2))
6 print('5. GridRegressionSVR : ',round(mean_squared_error(y_test, y_pred_grid),2))
7
8 print('-----')
9 print('MAE : \n1. SVR : ',round(mean_absolute_error(y_test, y_pred_svr),2))
10 print('2. Linear Regression : ',round(mean_absolute_error(y_test,lin_reg_pred),2))
11 print('3. Ridge Regression : ',round(mean_absolute_error(y_test, ridge_reg_pred),2))
12 print('4. Elastic Net Regression : ',round(mean_absolute_error(y_test, elastic_reg_pred),2))
13 print('5. GridRegressionSVR : ',round(mean_absolute_error(y_test, y_pred_grid),2))

-----
MSE :
1. SVR :  0.0
2. Linear Regression :  0.0
3. Ridge Regression :  0.0
4. Elastic Net Regression :  0.02
5. GridRegressionSVR :  0.0
-----
MAE :
1. SVR :  0.06
2. Linear Regression :  0.05
3. Ridge Regression :  0.05
4. Elastic Net Regression :  0.12
5. GridRegressionSVR :  0.05

In [76]: 1 print('Train Accuracy')
2 print('-----')
3 print('R-squared Accuracy : \n1.SVR : ',round(svr_r2_score*100,2))
4 print('2. Linear Regression : ',round(lin_reg_r2_score*100,2))
5 print('3. Ridge Regression : ',round(ridge_reg_r2_score*100,2))
6 print('4. Lasso Regression : ',round(lasso_reg_r2_score*100,2))
7 print('5. Elastic Net Regression : ',round(elastic_reg_r2_score*100,2))
8 print('6. GridSearchCV : ',round(grid_svr_r2_score*100,2))

Train Accuracy
-----
R-squared Accuracy :
1.SVR :  76.61
2. Linear Regression :  81.09
3. Ridge Regression :  81.11
4. Lasso Regression :  -1.32
5. Elastic Net Regression :  -1.32
6. GridSearchCV :  80.07

In [77]: 1 print('Test Accuracy')
2 print('-----')
3 print('Adjusted R-squared Accuracy : \n1.SVR : ',round(svr_r2_score*100,2))
4 print('2. Linear Regression : ',round(lin_reg_adj_r2_score*100,2))
5 print('3. Ridge Regression : ',round(ridge_reg_adj_r2_score*100,2))
6 print('4. Lasso Regression : ',round(lasso_reg_adj_r2_score*100,2))
7 print('5. Elastic Net Regression : ',round(elastic_reg_adj_r2_score*100,2))
8 print('6. GridSearchCV : ',round(grid_svr_adj_r2_score*100,2))

Test Accuracy
-----
Adjusted R-squared Accuracy :
1.SVR :  76.61
2. Linear Regression :  79.86
3. Ridge Regression :  79.88
4. Lasso Regression :  -7.91
5. Elastic Net Regression :  -7.91
6. GridSearchCV :  78.78
```

Observation :-

1	We can simply choose the best fit model according to this peformance score
---	--

In [ ]:	1	
---------	---	--