

Exploratory Data Analysis

```
In [1]: #Identification of variables and data types

In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

In [3]: data = pd.read_csv('Jamboree_Admission.csv')
data
Out[3]:
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
0	1	324	118	4	4.5	4.5	9.65	1	0.92
1	2	337	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	337	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 x 9 columns

```
In [4]: data.shape
Out[4]: (500, 9)
```

```
In [5]: data.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: 36.3 KB
```

```
In [6]: #Analysing the basic metrics

In [7]: data.describe(include="all")
Out[7]:
```

	Serial No.	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	500.000000	500.000000
mean	250.500000	316.472000	107.152000	3.114000	3.740000	3.484000	8.576440	0.560000	0.72174
std	144.481833	11.295148	6.081868	1.143512	0.991004	0.92545	0.604813	0.496884	0.14114
min	1.000000	290.000000	92.000000	1.000000	1.000000	1.00000	6.800000	0.000000	0.34000
25%	125.750000	308.000000	103.000000	2.000000	2.500000	3.00000	8.127500	0.000000	0.63000
50%	250.500000	317.000000	107.000000	3.000000	3.500000	3.50000	8.560000	1.000000	0.72000
75%	375.250000	325.000000	112.000000	4.000000	4.000000	4.00000	9.040000	1.000000	0.82000
max	500.000000	340.000000	120.000000	5.000000	5.000000	5.00000	9.920000	1.000000	0.97000

```
In [8]: data.describe(include=[np.number])
Out[8]:
```

	Serial No.	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	500.000000	500.000000
mean	250.500000	316.472000	107.152000	3.114000	3.740000	3.484000	8.576440	0.560000	0.72174
std	144.481833	11.295148	6.081868	1.143512	0.991004	0.92545	0.604813	0.496884	0.14114
min	1.000000	290.000000	92.000000	1.000000	1.000000	1.00000	6.800000	0.000000	0.34000
25%	125.750000	308.000000	103.000000	2.000000	2.500000	3.00000	8.127500	0.000000	0.63000
50%	250.500000	317.000000	107.000000	3.000000	3.500000	3.50000	8.560000	1.000000	0.72000
75%	375.250000	325.000000	112.000000	4.000000	4.000000	4.00000	9.040000	1.000000	0.82000
max	500.000000	340.000000	120.000000	5.000000	5.000000	5.00000	9.920000	1.000000	0.97000

```
In [9]: data.describe(include=[np.float64])
Out[9]:
```

	SOP	LOR	CGPA	Chance of Admit
count	500.000000	500.000000	500.000000	500.000000
mean	3.740000	3.484000	8.576440	0.72174
std	0.991004	0.92545	0.604813	0.14114
min	1.000000	1.000000	6.800000	0.340000
25%	2.500000	3.000000	8.127500	0.630000
50%	3.500000	3.500000	8.560000	0.720000
75%	4.000000	4.000000	9.040000	0.820000
max	5.000000	5.000000	9.920000	0.970000

```
In [10]: data.isnull().sum()/len(data)*100
Out[10]:
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
Serial No.	0.0								
GRE Score	0.0								
TOEFL Score	0.0								
University Rating	0.0								
SOP	0.0								
LOR	0.0								
CGPA	0.0								
Research	0.0								
Chance of Admit	0.0								
Dtype: float64									

```
In [11]: #Non-Graphical Analysis

In [12]: data['Research'].value_counts()/len(data)*100
Out[12]:
```

	Research
1	56.0
0	44.0
Name: Research, dtype: float64	

```
In [13]: groupby_research = data.groupby(['Research'])
In [14]: groupby_research.max()
Out[14]:
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
Research									
0	500	339	120	5	5.0	5.0	9.70	0.89	
1	498	340	120	5	5.0	5.0	9.92	0.97	

```
In [15]: data.groupby('GRE Score').agg(['Chance of Admit ':'count', 'mean', 'max', 'min'])
Out[15]:
```

	GRE Score	count	mean	max	min
GRE Score					
290	2	0.460000	0.47	0.45	
293	1	0.640000	0.64	0.64	
294	2	0.475000	0.49	0.46	
295	5	0.512000	0.69	0.37	
296	5	0.522000	0.61	0.44	
297	6	0.498333	0.59	0.34	
298	10	0.507000	0.69	0.34	
299	10	0.537000	0.68	0.38	
300	12	0.595933	0.71	0.36	
301	11	0.624545	0.68	0.44	
302	7	0.558571	0.65	0.46	
303	5	0.590000	0.77	0.36	
304	12	0.578333	0.71	0.38	
305	11	0.624545	0.71	0.53	
306	7	0.642857	0.73	0.48	
307	10	0.627000	0.79	0.47	
308	13	0.655385	0.77	0.46	
309	9	0.637778	0.76	0.48	
310	11	0.667273	0.76	0.54	
311	16	0.665000	0.79	0.42	
312	24	0.685417	0.81	0.50	
313	12	0.684167	0.79	0.53	
314	16	0.696250	0.84	0.54	
315	13	0.645385	0.79	0.39	
316	18	0.661667	0.77	0.49	
317	15	0.690000	0.84	0.49	
318	12	0.702500	0.78	0.63	
319	12	0.729167	0.82	0.65	
320	16	0.790000	0.82	0.64	
321	17	0.805294	0.93	0.69	
322	17	0.784706	0.92	0.61	
323	13	0.785385	0.92	0.45	
324	23	0.813913	0.92	0.70	
325	15	0.742667	0.91	0.52	
326	12	0.822500	0.91	0.74	
327	17	0.801176	0.93	0.61	
328	9	0.848889	0.94	0.78	
329	10	0.853000	0.90	0.72	
330	8	0.906250	0.93	0.85	
331	9	0.918889	0.94	0.86	
332	8	0.893750	0.94	0.79	
333	4	0.930000	0.96	0.89	
334	8	0.916250	0.97	0.78	
335	4	0.940000	0.96	0.78	
336	5	0.948000	0.97	0.92	
337	2	0.940000	0.97	0.92	
338	4	0.920000	0.95	0.91	
339	3	0.936667	0.96	0.89	
340	9	0.947778	0.97	0.90	

```
In [16]: data.groupby('TOEFL Score').agg(['Chance of Admit ':'count', 'mean', 'max', 'min'])
Out[16]:
```

	TOEFL Score	count	mean	max	min
TOEFL Score					
92	1	0.510000	0.51	0.51	
93	2	0.460000	0.46	0.46	
94	2	0.490000	0.56	0.42	
95	3	0.516667	0.62	0.44	
96	6	0.478657	0.54	0.34	
97	7	0.502857	0.65	0.38	
98	10	0.567000	0.76	0.34	
99	23	0.566957	0.76	0.36	
100	24	0.597083	0.78	0.42	
101	20	0.610500	0.79	0.38	
102	24	0.651667	0.80	0.50	
103	25	0.678000	0.76	0.52	
104	29	0.678966	0.82	0.42	
105	37	0.643108	0.79	0.39	
106	28	0.675429	0.81	0.52	
107	28	0.762814	0.94	0.53	
108	19	0.724211	0.87	0.45	
109	19	0.752632	0.85	0.62	
110	44	0.767045	0.92	0.46	
111	20	0.804000	0.93	0.52	
112	28	0.800000	0.96	0.55	
113	19	0.860526	0.96	0.76	
114	18	0.840556	0.96	0.61	
115	11	0.879091	0.95	0.70	
116	16	0.901875	0.96	0.78	
117	8	0.927000	0.96	0.87	
118	10	0.925000	0.94	0.90	
119	10	0.930000	0.97	0.89	
120	9	0.934444	0.97	0.86	

```
In [17]: data.groupby('Research').agg(['Chance of Admit ':'count', 'mean', 'max', 'min'])
Out[17]:
```

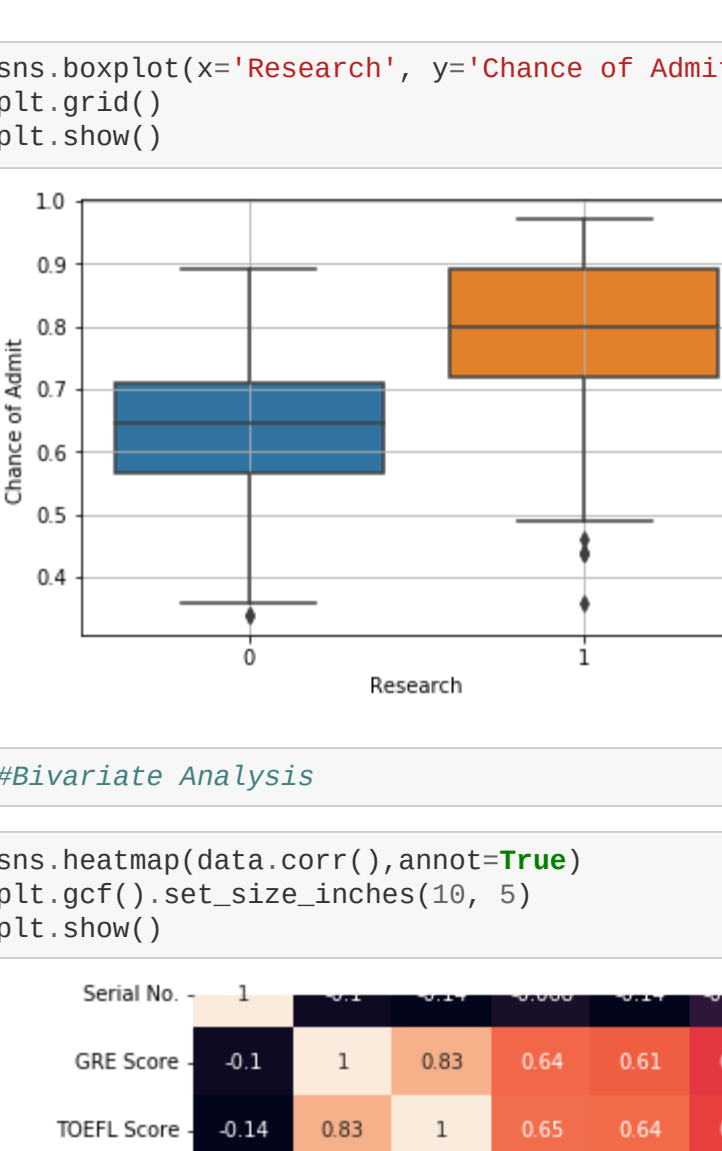
	Research	count	mean	max	min
Research					
0	220	0.634009	0.89	0.34	
1	280	0.789664	0.97	0.36	

```
In [18]: #Visual Analysis

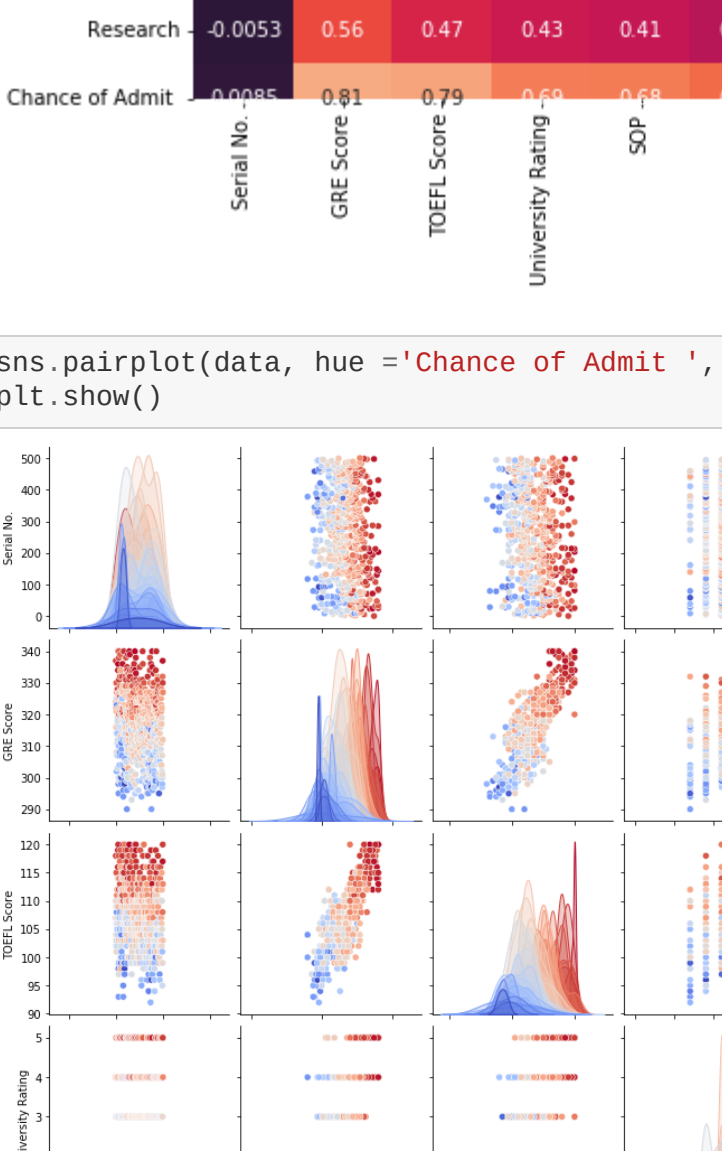
In [19]: data['GRE Score'].plot(kind='hist', bins=90)
Out[19]: <matplotlib.axes._subplots.AxesSubplot at 0x7f1de1771a58>
```



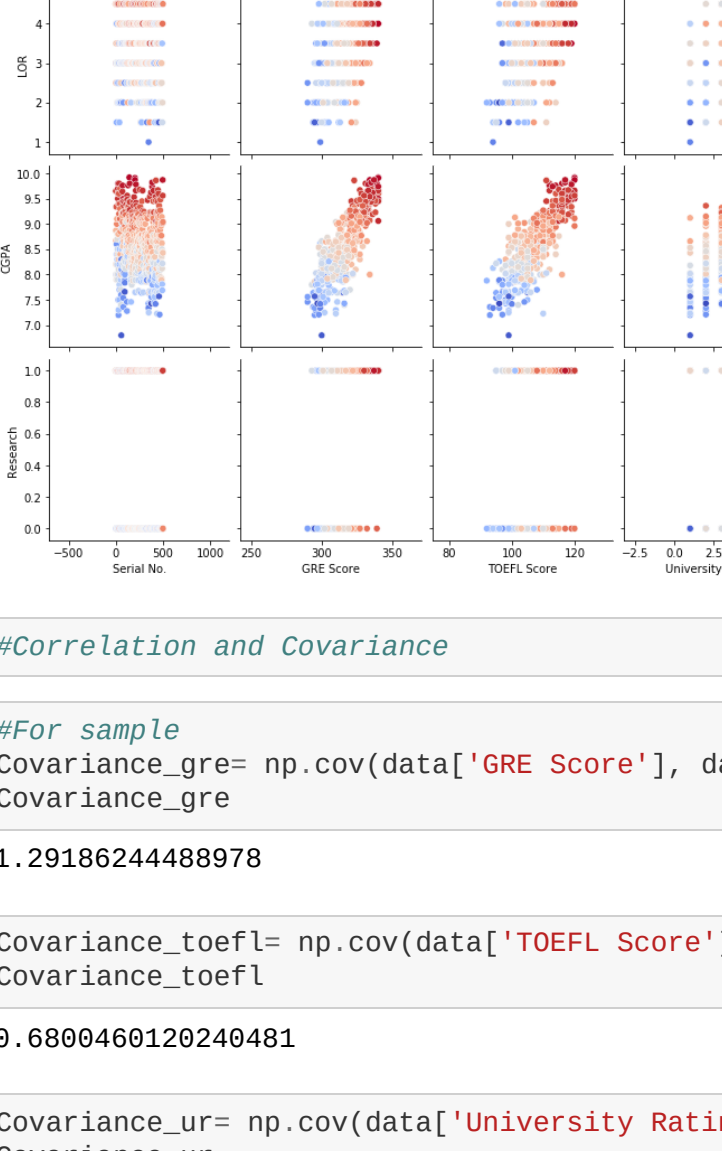
```
In [20]: data['TOEFL Score'].plot(kind='hist', bins=90)
Out[20]: <matplotlib.axes._subplots.AxesSubplot at 0x7f1dfe4ec5b0>
```



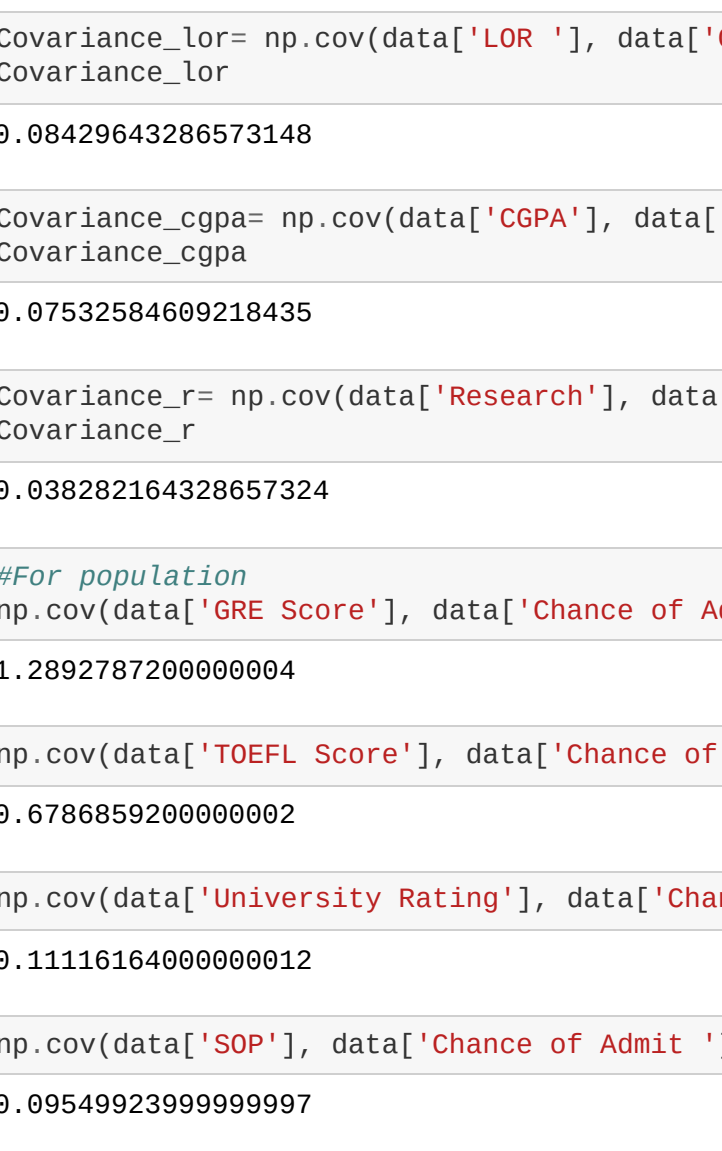
```
In [21]: data['CGPA'].plot(kind='hist', bins=90)
Out[21]: <matplotlib.axes._subplots.AxesSubplot at 0x7f1dfe4e0880>
```



```
In [22]: sns.boxplot(x='University Rating', y='Chance of Admit ', data=data)
plt.grid()
plt.show()
```



```
In [23]: sns.boxplot(x='Research', y='Chance of Admit ', data=data)
plt.grid()
plt.show()
```

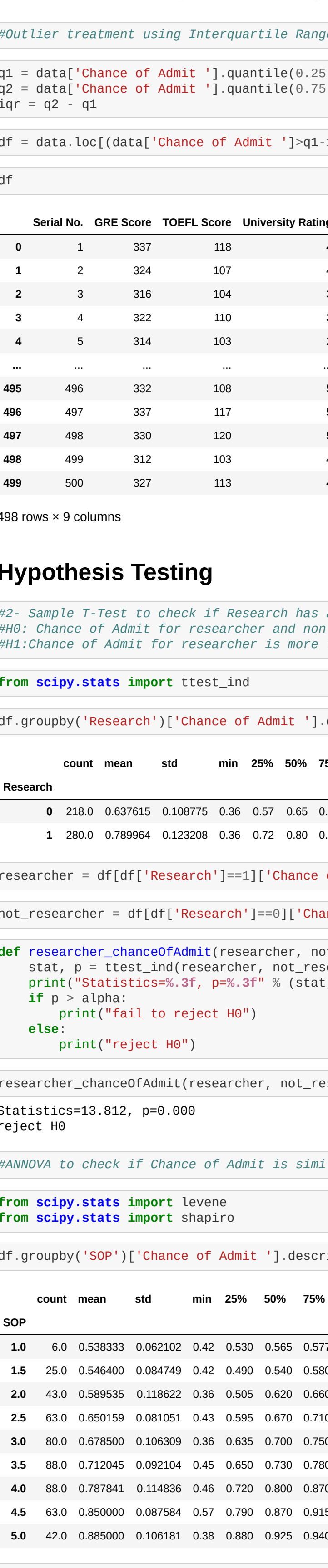


```
In [24]: #Bivariate Analysis

In [25]: sns.heatmap(data.corr(), annot=True)
plt.gcf().set_size_inches(10, 5)
plt.show()
```



```
In [26]: sns.pairplot(data, hue = 'Chance of Admit ', palette='coolwarm')
plt.show()
```



```
In [27]: #Correlation and Covariance

In [28]: #For sample
Covariance_gre= np.cov(data['GRE Score'], data['Chance of Admit '])[0][1]
Covariance_admit=1.2918624488978

In [29]: Covariance_toefl= np.cov(data['TOEFL Score'], data['Chance of Admit '])[0][1]
Covariance_toefl=0.68809460129249481

In [30]: Covariance_unr= np.cov(data['University Rating'], data['Chance of Admit '])[0][1]
Covariance_unr=0.11138440881763538

In [31]: Covariance_sop= np.cov(data['SOP'], data['Chance of Admit '])[0][1]
Covariance_sop=0.095606212448494

In [32]: Covariance_lor= np.cov(data['LOR'], data['Chance of Admit '])[0][1]
Covariance_lor=0.08429643286573148

In [33]: Covariance_cgpa= np.cov(data['CGPA'], data['Chance of Admit '])[0][1]
Covariance_cgpa=0.07532584609218435

In [34]: Covariance_re= np.cov(data['Research'], data['Chance of Admit '])[0][1]
Covariance_re=0.038282164328657324

In [35]: #For population
Covariance_gre= np.cov(data['GRE Score'], data['Chance of Admit '], bias=True)[0][1]
Covariance_gre=1.2892787208080804

In [36]: np.cov(data['TOEFL Score'], data['Chance of Admit '], bias=True)[0][1]
Covariance_toefl=0.6786859208080802

In [37]: np.cov(data['University Rating'], data['Chance of Admit '], bias=True)[0][1]
Covariance_unr=0.1116164008080812

In [38]: np.cov(data['SOP'], data['Chance of Admit '], bias=True)[0][1]
Covariance_sop=0.0954992399999997

In [39]: np.cov(data['LOR'], data['Chance of Admit '], bias=True)[0][1]
Covariance_lor=0.084127840000000802

In [40]: np.cov(data['CGPA'], data['Chance of Admit '], bias=True)[0][1]
Covariance_cgpa=0.07517519439999988

In [41]: np.cov(data['Research'], data['Chance of Admit '], bias=True)[0][1]
Covariance_re=0.038205600000000801

In [42]: #Correlation Coefficient
np.corrcoef(data['GRE Score'], data['Chance of Admit '])[0][1]
Covariance_gre=0.8103566354632607

In [43]: np.corrcoef(data['TOEFL Score'], data['Chance of Admit '])[0][1]
Covariance_toefl=0.7922276143050841

In [44]: np.corrcoef(data['University Rating'], data['Chance of Admit '])[0][1]
Covariance_unr=0.6991323687886996

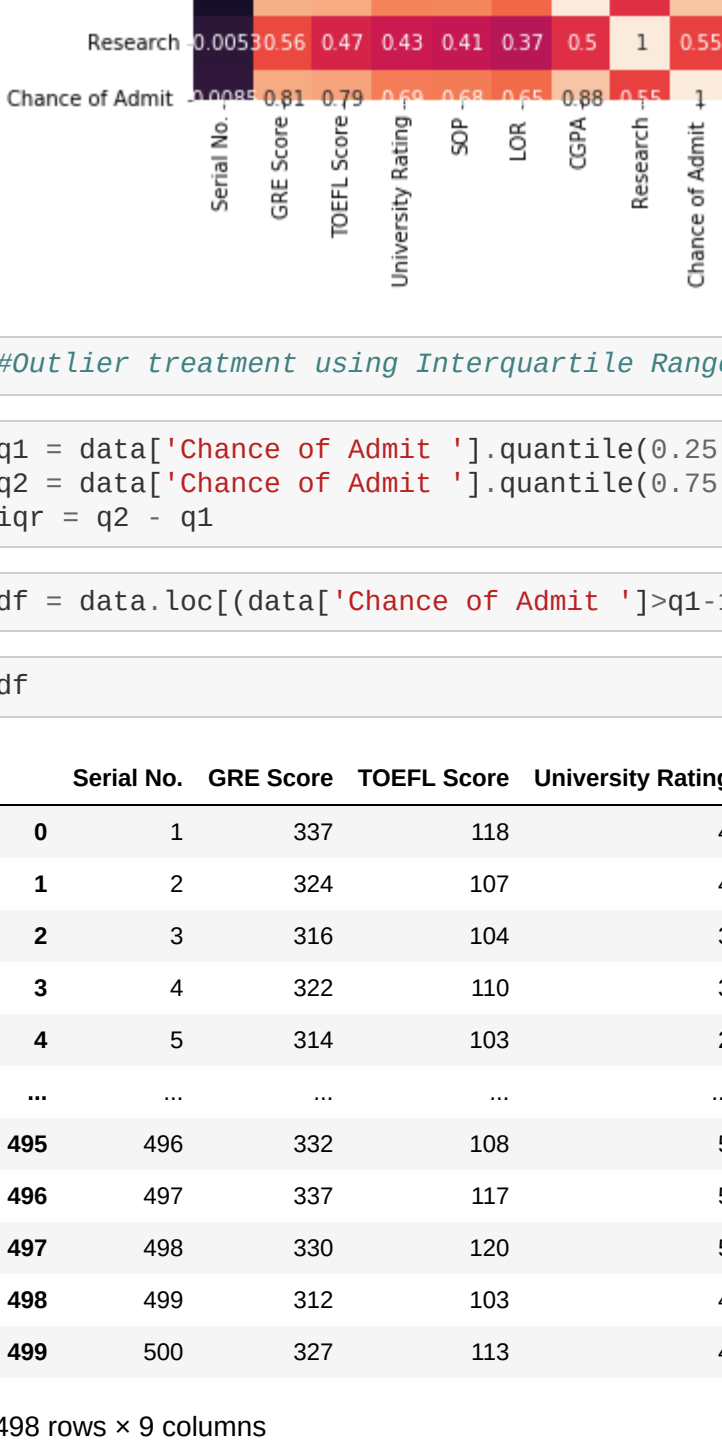
In [45]: np.corrcoef(data['SOP'], data['Chance of Admit '])[0][1]
Covariance_sop=0.6841365241316723

In [46]: np.corrcoef(data['LOR'], data['Chance of Admit '])[0][1]
Covariance_lor=0.6453645135280114

In [47]: np.corrcoef(data['CGPA'], data['Chance of Admit '])[0][1]
Covariance_cgpa=0.8824125749045737

In [48]: np.corrcoef(data['Research'], data['Chance of Admit '])[0][1]
Covariance_re=0.5458710294711387

In [49]: #Visualization
sns.heatmap(data.corr(), annot=True)
Out[49]: <matplotlib.axes._subplots.AxesSubplot at 0x7f1dd83f8c18>
```



```
In [50]: #Outlier treatment using Interquartile Range(IQR)

In [51]: q1 = data['Chance of Admit '].quantile(0.25)
q2 = data['Chance of Admit '].quantile(0.75)
IQR = q2 - q1

In [52]: df = data.loc[(data['Chance of Admit ']>q1-1.5*IQR) & (data['Chance of Admit ']<=q2+1.5*IQR)]

In [53]: df
Out[53]:
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
0	1	324	118	4	4.5	4.5	9.65	1	0.92
1	2	337	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						


```
In [75]: from scipy.stats import chi2_contingency

stat, p, dof, expected = chi2_contingency(gre.coa)

# interpret p-value
alpha = 0.05
print("p value is " + str(p))
if p <= alpha:
    print('Dependent (reject H0)')
else:
    print('Independent (H0 holds true)')

p value is 1.0
Independent (H0 holds true)
```

```
In [76]: from scipy.stats import chi2_contingency

stat, p, dof, expected = chi2_contingency(tofel.coa)

# interpret p-value
alpha = 0.05
print("p value is " + str(p))
if p <= alpha:
    print('Dependent (reject H0)')
else:
    print('Independent (H0 holds true)')

p value is 1.0
Independent (H0 holds true)
```

Linear Regression

```
In [77]: X = df.drop(['Chance of Admit '], axis=1)
         Y = df['Chance of Admit ']
```

```
In [78]: from sklearn.model_selection import train_test_split
```

```
In [79]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2)
```

```
In [80]: from sklearn import linear_model
         from sklearn.metrics import mean_squared_error, r2_score
```

```
In [81]: model = linear_model.LinearRegression()
         model.fit(X_train, Y_train)
```

```
Out[81]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None,
                          normalize=False)
```

```
In [82]: Y_pred_train = model.predict(X_train)
```

```
In [83]: print('Coefficients:', model.coef_)
         print('Intercept:', model.intercept_)
         print('Mean squared error (MSE): %.2f' % mean_squared_error(Y_train, Y_pred_train))
         print('Coefficient of determination (R^2): %.2f' % r2_score(Y_train, Y_pred_train))

Coefficients: [9.46356118e-05 1.78528194e-03 2.89878249e-03 2.97587888e-03
              7.86984587e-03 1.31683648e-02 1.28375548e-01 2.44657482e-02]
Intercept: -1.3847887473693473
Mean squared error (MSE): 0.00
Coefficient of determination (R^2): 0.84
```

```
In [84]: Y_pred_test = model.predict(X_test)
```

```
In [85]: print('Coefficients:', model.coef_)
         print('Intercept:', model.intercept_)
         print('Mean squared error (MSE): %.2f'
               % mean_squared_error(Y_test, Y_pred_test))
         print('Coefficient of determination (R^2): %.2f'
               % r2_score(Y_test, Y_pred_test))

Coefficients: [9.46356118e-05 1.78528194e-03 2.89878249e-03 2.97587888e-03
              7.86984587e-03 1.31683648e-02 1.28375548e-01 2.44657482e-02]
Intercept: -1.3847887473693473
Mean squared error (MSE): 0.00
Coefficient of determination (R^2): 0.80
```

```
In [86]: yintercept = '%.2f' % model.intercept_
         LogP = '%.2f LogP' % model.coef_[0]
         MW = '%.4f MW' % model.coef_[1]
         RB = '%.4f RB' % model.coef_[2]
         AP = '%.2f AP' % model.coef_[3]
```

```
In [87]: print('LogS = ' + ' ' + yintercept + ' ' + LogP + ' ' + MW + ' ' + RB + ' ' + AP)

LogS = -1.38 0.00 LogP 0.0018 MW + 0.0029 RB 0.00 AP
```

```
In [88]: #Vertical plot
```

```
In [89]: plt.figure(figsize=(5,11))

plt.subplot(2, 1, 1)
plt.scatter(x=Y_train, y=Y_pred_train, c="#7CAE00", alpha=0.3)

z = np.polyfit(Y_train, Y_pred_train, 1)
p = np.polyval(z)
plt.plot(Y_test, p(Y_test), "#F8766D")

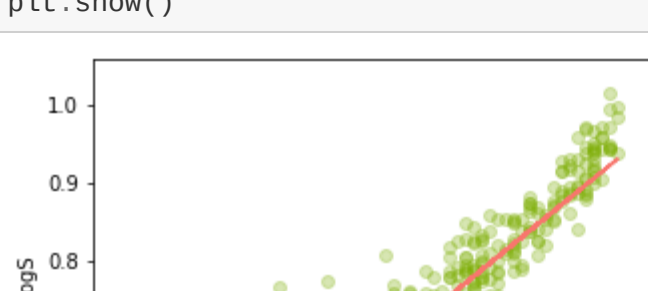
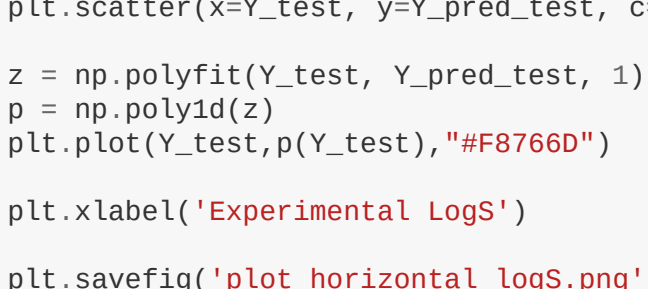
plt.ylabel('Predicted LogS')

# 2 row, 1 column, plot 2
plt.subplot(2, 1, 2)
plt.scatter(x=Y_test, y=Y_pred_test, c="#619CFF", alpha=0.3)

z = np.polyfit(Y_test, Y_pred_test, 1)
p = np.polyval(z)
plt.plot(Y_test, p(Y_test), "#F8766D")

plt.ylabel('Predicted LogS')
plt.xlabel('Experimental LogS')

plt.savefig('plot_vertical_LogS.png')
plt.savefig('plot_vertical_LogS.pdf')
plt.show()
```



```
In [90]: #Horizontal plot
```

```
In [91]: plt.figure(figsize=(11,5))

plt.subplot(1, 2, 1)
plt.scatter(x=Y_train, y=Y_pred_train, c="#7CAE00", alpha=0.3)

z = np.polyfit(Y_train, Y_pred_train, 1)
p = np.polyval(z)
plt.plot(Y_test, p(Y_test), "#F8766D")

plt.ylabel('Predicted LogS')
plt.xlabel('Experimental LogS')

# 1 row, 2 column, plot 2
plt.subplot(1, 2, 2)
plt.scatter(x=Y_test, y=Y_pred_test, c="#619CFF", alpha=0.3)

z = np.polyfit(Y_test, Y_pred_test, 1)
p = np.polyval(z)
plt.plot(Y_test, p(Y_test), "#F8766D")

plt.xlabel('Experimental LogS')

plt.savefig('plot_horizontal_LogS.png')
plt.savefig('plot_horizontal_LogS.pdf')
plt.show()
```

