

## IMPORT LIBRARIES AND DATASET

```
In [152]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt

# Read the csv file
admission_df = pd.read_csv('Admission_Predict.csv')
admission_df.head()
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
0	337	106	118	4	4.5	4.5	9.65	1	0.92
1	324	107	107	4	4.0	4.5	8.67	1	0.76
2	316	104	104	3	3.0	3.5	8.00	1	0.72
3	322	110	110	3	3.5	2.5	8.67	1	0.80
4	314	103	103	2	2.0	3.0	8.21	0	0.65

```
In [153]: # Drop Serial No. column
admission_df.drop('Serial No.', axis = 1, inplace = True)
```

```
In [154]: admission_df
```

	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	8.67	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
...	...	...	...	...	...	...	...	...
495	332	108	5	4.5	4.0	9.02	1	0.87
496	337	117	5	5.0	5.0	9.87	1	0.96
497	330	120	5	4.5	5.0	9.56	1	0.93
498	312	103	4	4.0	5.0	8.43	0	0.73
499	327	113	4	4.5	4.5	9.04	0	0.84

500 rows x 9 columns

## PERFORM EXPLORATORY DATA ANALYSIS

```
In [155]: # Check for null values
admission_df.isnull().sum()
```

```
In [156]: # Check for dataframe information
admission_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 9 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
 8   dtype: float64(4), int64(4)
memory usage: 31.4 KB
```

```
In [157]: # Statistical summary of the dataframe
admission_df.describe()
```

	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
mean	316.472000	107.192000	3.114000	3.374000	3.484000	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	1.000000	6.800000	0.000000	0.34000
25%	308.000000	103.000000	2.000000	2.500000	3.000000	8.127500	0.000000	0.63000
50%	317.000000	107.000000	3.000000	3.500000	3.500000	8.560000	1.000000	0.72000
75%	325.000000	112.000000	4.000000	4.000000	4.000000	9.040000	1.000000	0.82000
max	340.000000	120.000000	5.000000	5.000000	5.000000	9.920000	1.000000	0.97000

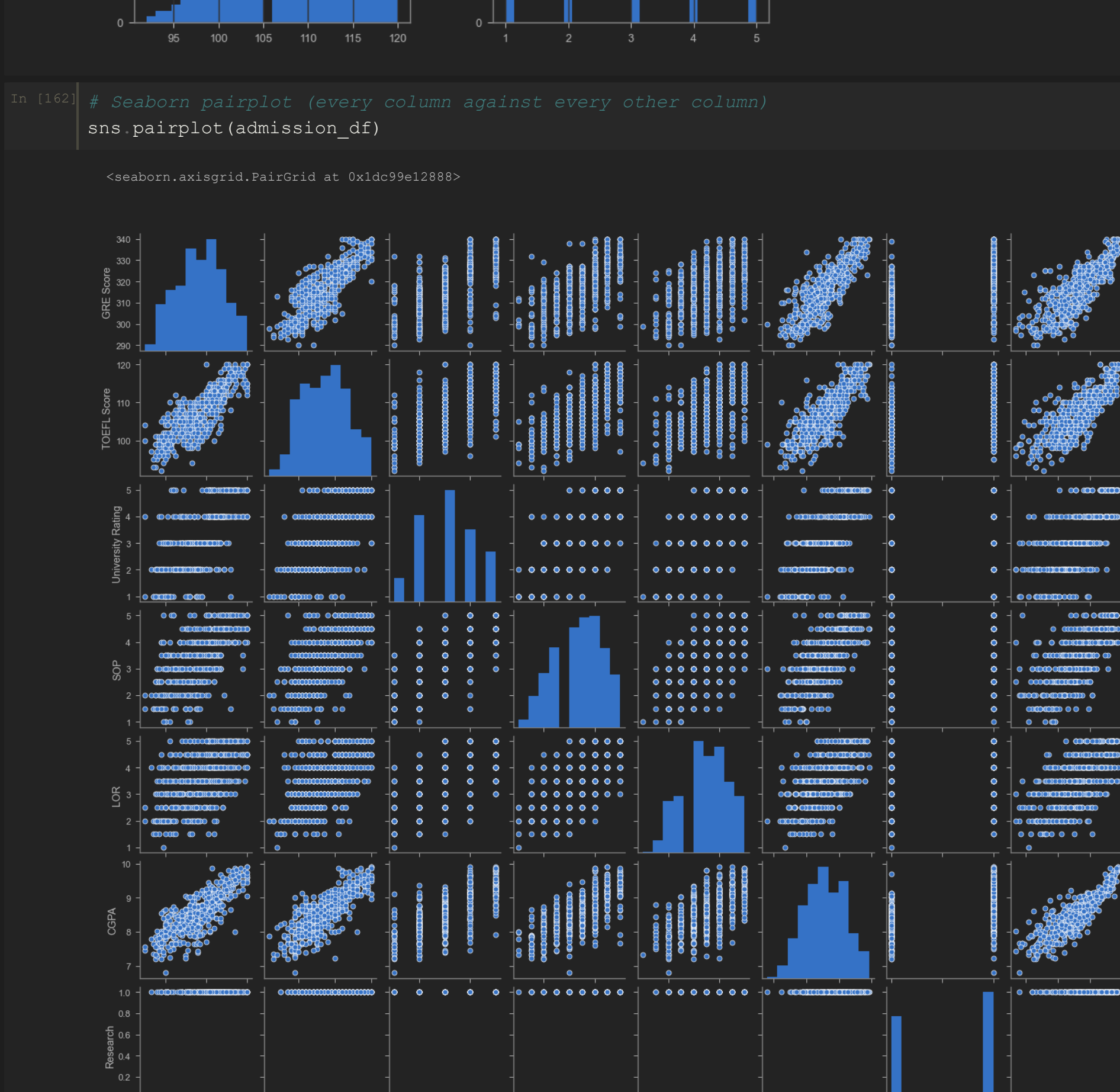
```
In [158]: # Grouping by university ranking
df_university = admission_df.groupby(by = 'University Rating').mean()
```

```
In [159]: df_university
```

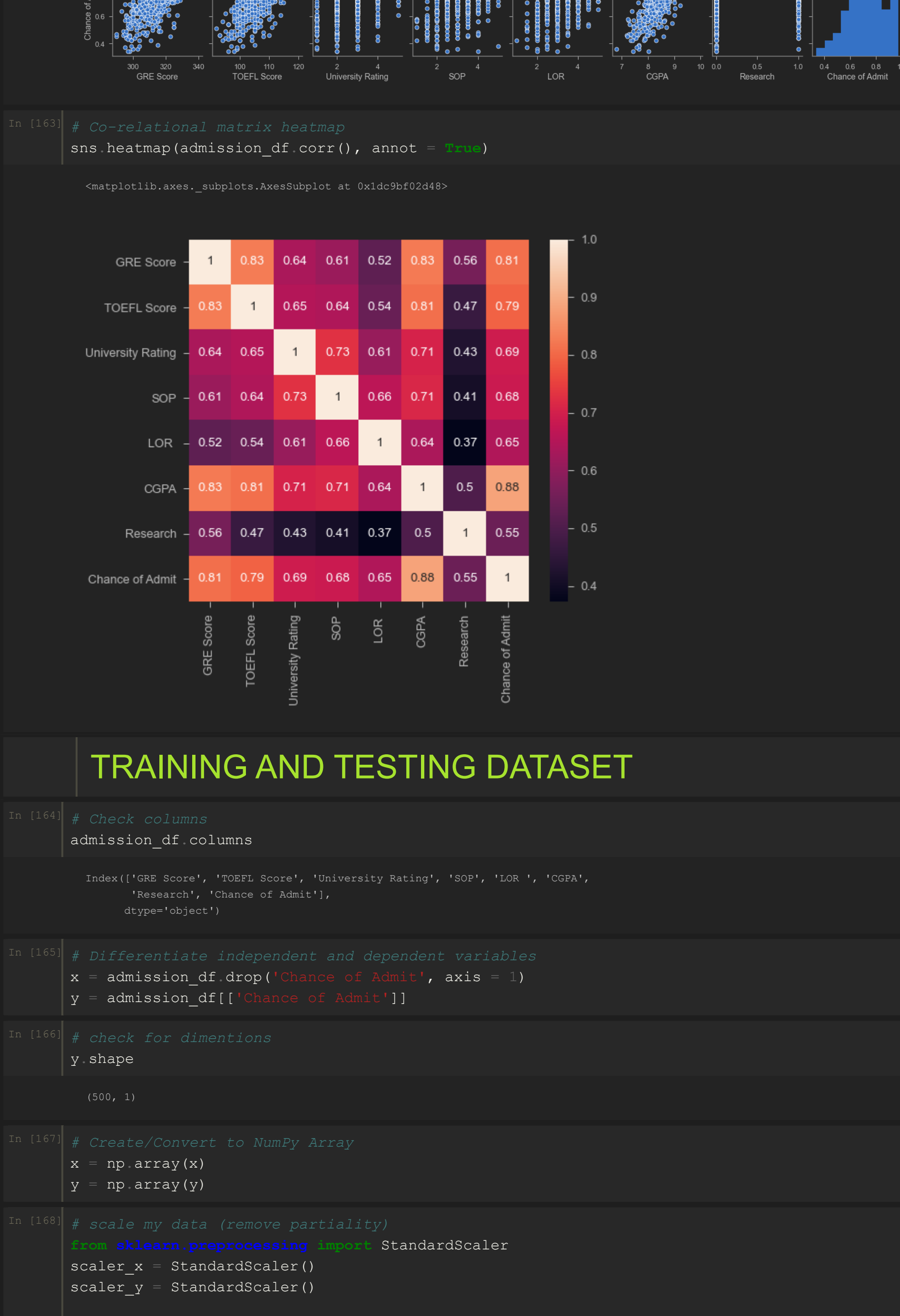
	GRE Score	TOEFL Score	SOP	LOR	CGPA	Research	Chance of Admit
University Rating							
1	304.911765	100.205882	1.941176	2.42471	7.798529	0.294118	0.562059
2	309.114321	103.444444	2.682540	2.956349	8.177778	0.293651	0.626111
3	315.000864	106.314815	3.308642	3.401235	8.500123	0.537037	0.702901
4	323.804762	110.961905	4.000000	3.747619	8.936667	0.769592	0.801619
5	327.880411	113.438556	4.479452	4.404110	9.278082	0.876712	0.888032

## PERFORM DATA VISUALIZATION

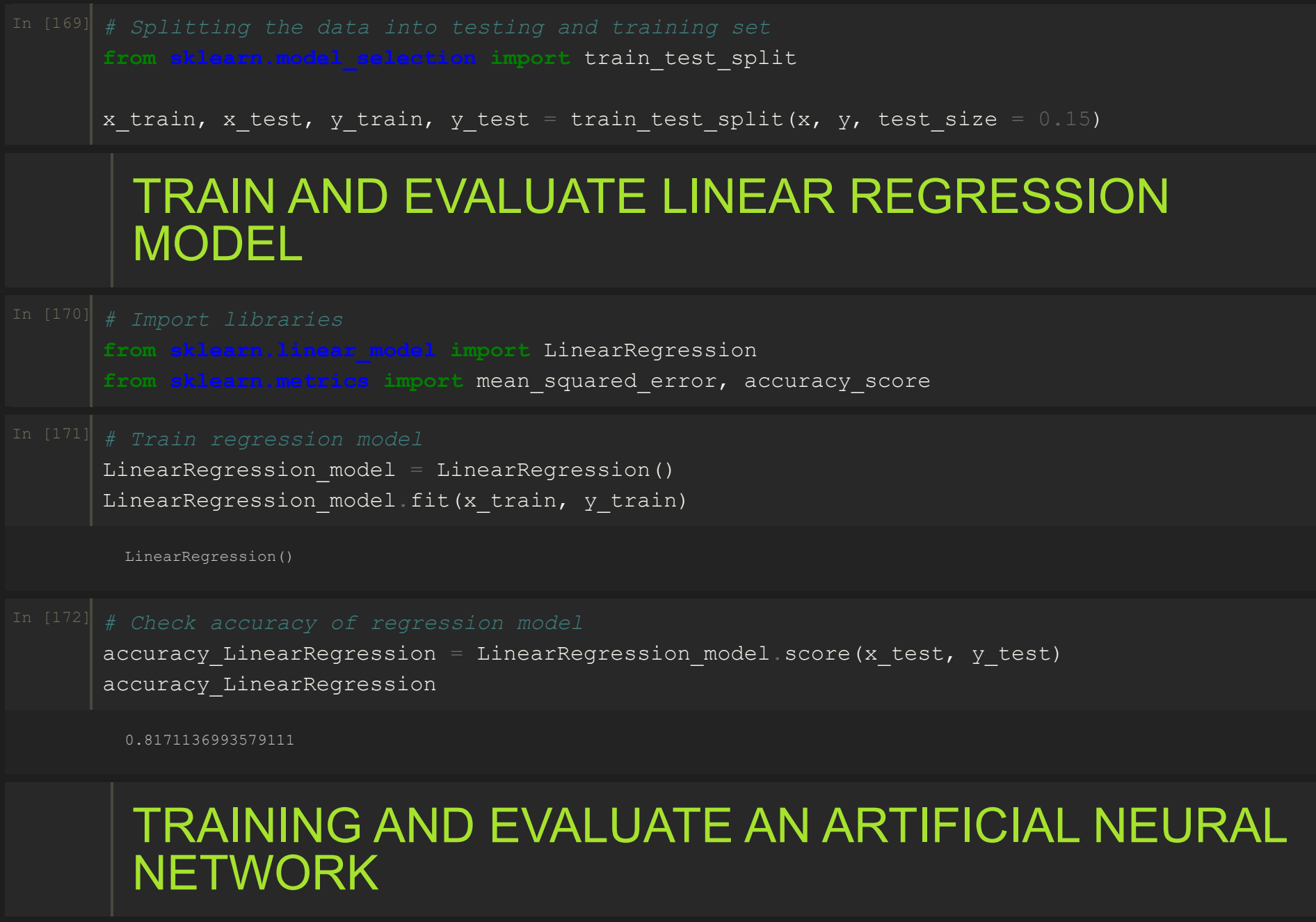
```
In [160]: # Histogram of every column
admission_df.hist(bins = 20, figsize = (20, 20))
```



```
In [161]: # Seaborn pairplot (every column against every other column)
sns.pairplot(admission_df)
```



```
In [162]: # Correlational matrix heatmap
sns.heatmap(admission_df.corr(), annot = True)
```



## TRAINING AND TESTING DATASET

```
In [164]: # Check columns
Index(['GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR', 'CGPA', 'Research', 'Chance of Admit'],
      dtype='object')
```

```
In [165]: # Differentiate Independent and Dependent variables
x = admission_df.drop('Chance of Admit', axis = 1)
y = admission_df[['Chance of Admit']]
```

```
In [166]: # Check for dimensions
y.shape
```

```
(500, 1)
```

```
In [167]: # Create/Convert to NumPy Array
x = np.array(x)
y = np.array(y)
```

```
In [168]: # Scale my data (remove partiality)
from sklearn.preprocessing import StandardScaler
scaler_x = StandardScaler()
scaler_y = StandardScaler()
```

```
x = scaler_x.fit_transform(x)
y = scaler_y.fit_transform(y)
```

```
In [169]: # Splitting the data into testing and training set
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.15)
```

## TRAIN AND EVALUATE LINEAR REGRESSION MODEL

```
In [170]: # Import libraries
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, accuracy_score
```

```
In [171]: # Train regression model
LinearRegression model = LinearRegression()
LinearRegression model.fit(x_train, y_train)
LinearRegression()
```

```
In [172]: # Check accuracy of regression model
accuracy_linearRegression = LinearRegression model.score(x_test, y_test)
accuracy_linearRegression
```

```
0.8171136931979111
```

## TRAINING AND EVALUATE AN ARTIFICIAL NEURAL NETWORK

```
In [173]: # Import libraries
import tensorflow as tf
from tensorflow.keras.layers import Dense, Activation, Dropout
from tensorflow.keras.optimizers import Adam
```

```
In [174]: # Building Artificial Neural Network
ANN_model = keras.Sequential()
```

```
ANN_model.add(Dense(10, input_dim = 7))
ANN_model.add(Activation('relu'))
```

```
ANN_model.add(Dense(10))
ANN_model.add(Activation('relu'))
```

```
ANN_model.add(Dense(150))
ANN_model.add(Activation('relu'))
```

```
ANN_model.add(Dense(150))
ANN_model.add(Activation('relu'))
```

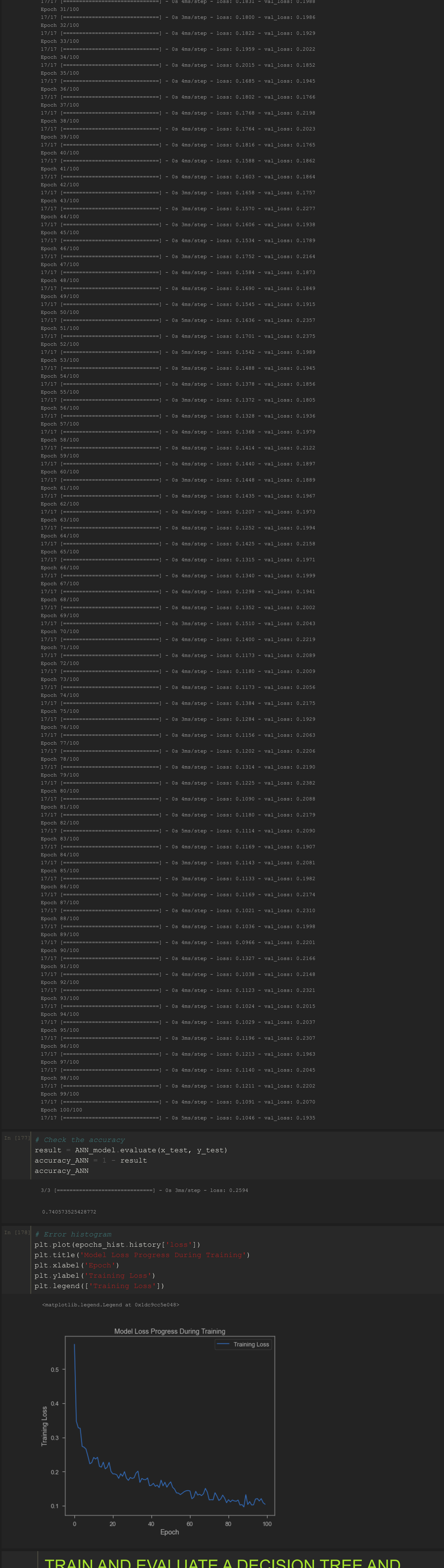
```
ANN_model.add(Dense(50))
ANN_model.add(Activation('linear'))
```

```
ANN_model.compile(loss = 'mse', optimizer = 'adam')
ANN_model.summary()
```

```
Model: "sequential_3"
Layer (type)                 Output Shape         Param #
-----
dense_11 (Dense)              (None, 50)           400
activation_2 (Activation)      (None, 50)           0
dense_12 (Dense)              (None, 150)          7650
activation_3 (Activation)      (None, 150)          0
dropout_4 (Dropout)           (None, 150)          0
dense_13 (Dense)              (None, 150)          22650
activation_4 (Activation)      (None, 150)          0
dropout_5 (Dropout)           (None, 150)          0
dense_14 (Dense)              (None, 50)           7350
activation_5 (Activation)      (None, 50)           0
dense_15 (Dense)              (None, 1)             51
Total params: 36,301
Trainable params: 36,301
Non-trainable params: 0
```

```
In [175]: # Compile the model
ANN_model.compile(optimizer = 'adam', loss = 'mean_squared_error')
```

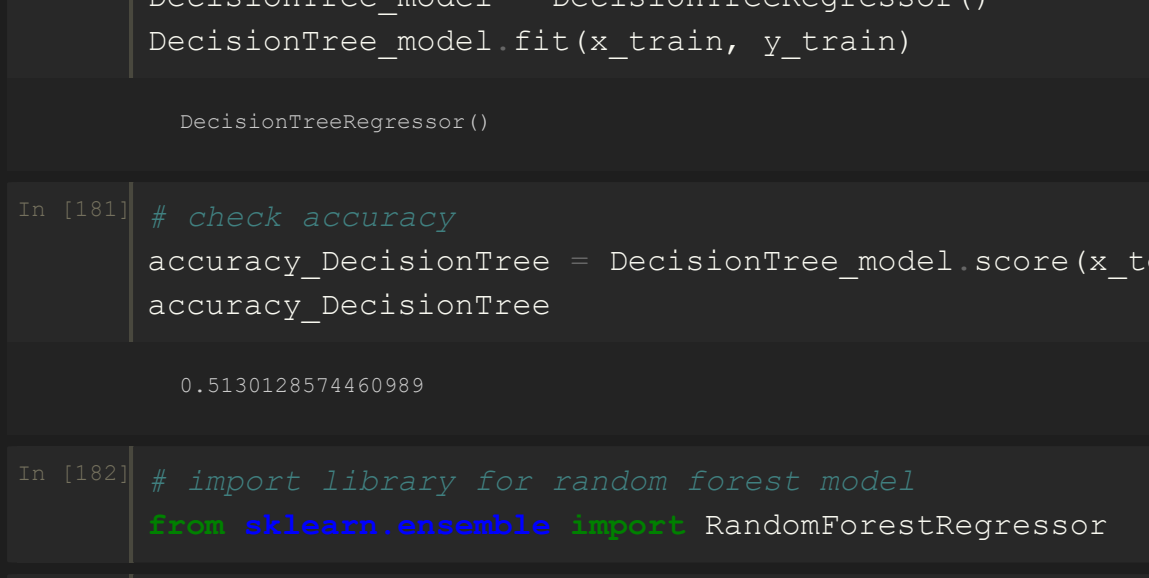
```
In [176]: # Train the network
epochs_hist = ANN_model.fit(x_train, y_train, epochs = 100, batch_size = 20, validation_split = 0.2)
```



```
In [177]: # Check the accuracy
result = ANN_model.evaluate(x_test, y_test)
accuracy ANN = 1 - result
accuracy ANN
```

```
0.740973525428772
```

```
In [178]: # Error histogram
plt.plot(epochs_hist.history['loss'])
plt.title('Model Loss Progress During Training')
plt.xlabel('Epoch')
plt.ylabel('Training Loss')
plt.legend(['Training Loss'])
```



## TRAIN AND EVALUATE A DECISION TREE AND RANDOM FOREST MODELS

```
In [179]: # Import library for decision tree model
from sklearn.tree import DecisionTreeRegressor
```

```
In [180]: # Train decision tree model
DecisionTree_model = DecisionTreeRegressor()
DecisionTree_model.fit(x_train, y_train)
DecisionTreeRegressor()
```

```
In [181]: # Check accuracy
accuracy_DecisionTree = DecisionTree_model.score(x_test, y_test)
accuracy_DecisionTree
```

```
0.5130128571460989
```

```
In [182]: # Import library for random forest model
from sklearn.ensemble import RandomForestRegressor
```

```
In [183]: # Train random forest model
RandomForest_model = RandomForestRegressor(n_estimators = 80, max_depth = 10)
RandomForest_model.fit(x_train, y_train)
```

```
C:\Users\Tejas\anaconda3\lib\site-packages\ipykernel_launcher.py:1: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, 1), for example using ravel().
This is separate from the ipykernel package so we can avoid doing imports until
RandomForestRegressor(max_depth=10, n_estimators=80)
```

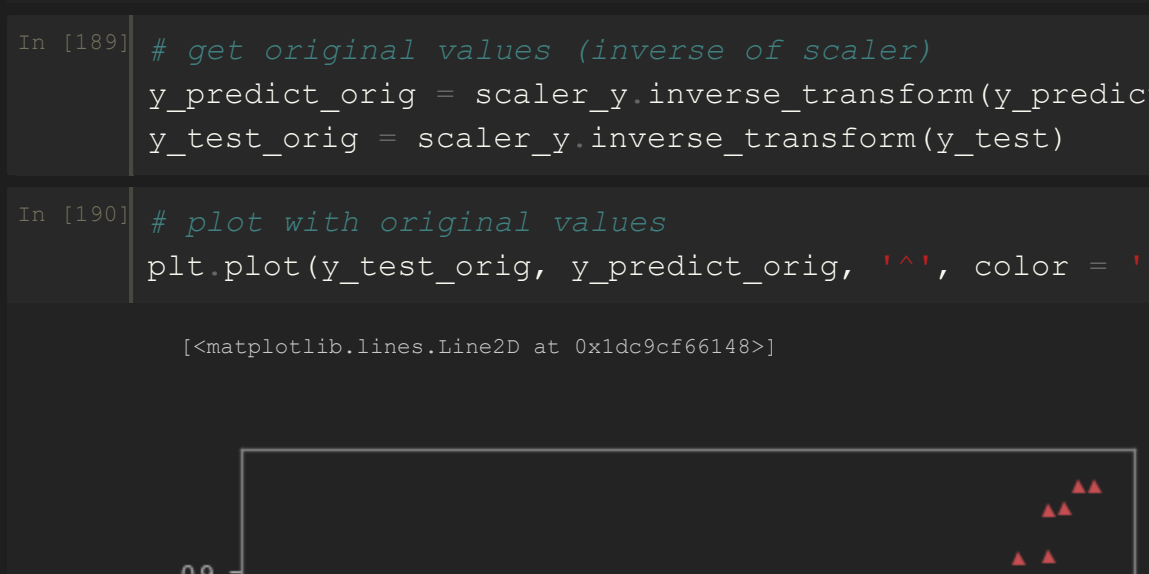
```
In [184]: # Check accuracy
accuracy_RandomForest = RandomForest_model.score(x_test, y_test)
accuracy_RandomForest
```

```
0.77208393790997
```

## CALCULATE REGRESSION MODEL KPIS

```
In [185]: # get predictions from test set
y_predict = LinearRegression_model.predict(x_test)
```

```
In [186]: # plot y_test vs y_predict
plt.plot(y_test, y_predict, 'r', color = 'r')
```



```
In [188]: # get original values (inverse of scaler)
y_predict_orig = scaler_y.inverse_transform(y_predict)
y_test_orig = scaler_y.inverse_transform(y_test)
```

```
In [189]: # plot with original values
plt.plot(y_test_orig, y_predict_orig, 'r', color = 'r')
```



```
In [190]: # get the length of test data (n) and no. of independent variables (k)
k = x_test.shape[1]
n = len(x_test)
n
```

```
75
```



```
In [195] # import libraries
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from math import sqrt

In [203] # calculate RMSE
RMSE = sqrt(mean_squared_error(y_test_orig, y_predict_orig))

In [197] # calculate MSE
MSE = mean_squared_error(y_test_orig, y_predict_orig)

In [198] # calculate MAE
MAE = mean_absolute_error(y_test_orig, y_predict_orig)

In [199] # calculate R2 score
r2 = r2_score(y_test_orig, y_predict_orig)

In [200] # calculate adjusted R2 score
adj_r2 = 1 - ((1 - r2)*(n - 1) / (n - k - 1))

In [204] # print errors
print('RMSE =', RMSE, '\nMSE =', MSE, '\nMAE =', MAE, '\nR square =', r2, '\nAdjusted R square =', adj_r2)

RMSE = 0.05897407078888444
MSE = 0.003480923184377506
MAE = 0.04359252597839725
R square = 0.817212693579109
Adjusted R square = 0.798064154102239

In [209] # predict with custom inputs
new_predict = LinearRegression_model.predict(np.array([[340, 120, 1, 5.0, 5.0, 10.0, 1]]))

new_predict

array([[71.33029724)])
```

THE END