```
In [ ]:     1  #*******************************************************************************
            2  #
            3  #              Project 2 : Prediction of Overall Rating of the Soccer Player
            4  #
            5  #
            6  #*******************************************************************************
            7
            8
            9  # Step0  : importing Libraries
           10
           11  import numpy as np
           12  import pandas as pd
           13  import sqlite3
           14
           15  from sklearn.tree import DecisionTreeRegressor
           16  from sklearn.linear_model import LinearRegression
           17  from sklearn.model_selection import train_test_split
           18  from sklearn.metrics import mean_squared_error,r2_score
           19  from sklearn.preprocessing import LabelEncoder
           20
           21  import seaborn as sns
           22  import matplotlib.pyplot as plt
           23  %matplotlib inline
           24  from math import sqrt
           25
           26  # Step1  : Loading Data from SQl DataBase
           27
           28  # Creating a connection to the sqlite DataBase and fetching Data
           29
           30  cnx = sqlite3.connect('database.sqlite')
           31  df = pd.read_sql_query("SELECT * FROM Player_Attributes", cnx)
           32
           33
           34  # Data from Player_Attributes Table
           35  print(" The Data from the Player Attributes Table giving information of Soccer Player")
           36  print('-'*80)
           37  print(df.head(2))
           38
           39  print("The attributes of the Player to be compared with are :\n ",'-'*80)
           40  print(df.columns)
           41
```

```
In [ ]:    1   # Step2 : Data Preprocessing
           2
           3   print("The number of rows and columns\n")
           4   print('-'*80)
           5   print(df.shape)
           6
           7   # Data  has 183978 rows and 42 columns
           8
           9   # Understanding the data
          10   print('-'*80)
          11   print(df.describe())
          12   # all ids and Date columns have high std  , so we can drop them
          13
          14   print('-'*80)
          15   print(df.info())
          16
          17
          18
          19   # check for if any columns have nulls
          20   print(" Checking for nulls in any of the features \n",'-'*80)
          21
          22   print('-'*80)
          23   print((df.isnull().sum()/df.shape[0])*100)
          24
          25   # since the percentage of nulls is less than total number of rows , we can drop the nulls
          26   df = df.dropna()
          27
          28   print("After Dropping the nulls check for the count of nulls and the shape \n",'-'*80)
          29
          30   print((df.isnull().sum()/df.shape[0])*100)
          31   print('-'*80)
          32
          33   print("Number of rows after dropping the nulls \n",'-'*80)
          34   print(df.shape)
          35
          36
```

In [ ]:

```python
#Step 3  :  Data Analysing

# From the Data we see there are Categorical Fields like attacking work rate , Defensive Work Rate
#'preferred_foot', 'attacking_work_rate','defensive_work_rate'
# finding unique values for each
print(df.preferred_foot.unique())
print(df.attacking_work_rate.unique())
print(df.defensive_work_rate.unique())

print(" Converting the Categorical attributes to numerical values \n",'-'*80)

number = LabelEncoder()
df['preferred_foot'] = number.fit_transform(df['preferred_foot'].astype('str'))
df['attacking_work_rate'] = number.fit_transform(df['attacking_work_rate'].astype('str'))
df['defensive_work_rate'] = number.fit_transform(df['defensive_work_rate'].astype('str'))

print(df.preferred_foot.unique())
print(df.attacking_work_rate.unique())
print(df.defensive_work_rate.unique())

# Dropping the ID fields
df.drop(['id','player_fifa_api_id','player_api_id', 'date'], axis=1,inplace=True)

print(" The Data after dropping the ids \n",'-'*80)
print(df.head())

```

```
In [ ]:    1  #Step 4 : Visualising the Data by using pairplot with respect to overall rating
           2
           3  df.head()
           4
           5  features = df.columns
           6
           7  target = 'overall_rating'
           8  count = len(df.columns)
           9
          10  # as lots of records , pairplot done only for sample size
          11
          12  i = 1
          13  while i < count:
          14      sns.pairplot(df.sample(1000,random_state=2),x_vars= features[i:i+5] ,y_vars = target, aspect=1,kind='reg')
          15      i = i + 5
          16
          17
```

```python
# Step 5 : Splitting the Data into Training and Test Set

# from the Visualisation we see that colummns
# 'preferred_foot', 'attacking_work_rate','defensive_work_rate' has no much effect on the overall rating
# so we can ignore them
# also by checking for fields 'sliding_tackle', 'gk_handling','gk_reflexes' not much effect on accuracy
# by inclusion, so dropping it

features = ['potential', 'crossing', 'finishing', 'heading_accuracy',
        'short_passing', 'volleys', 'dribbling', 'curve', 'free_kick_accuracy',
        'long_passing', 'ball_control', 'acceleration', 'sprint_speed',
        'agility', 'reactions', 'balance', 'shot_power', 'jumping', 'stamina',
        'strength', 'long_shots', 'aggression','interceptions', 'positioning',
         'vision', 'penalties', 'marking' ,'standing_tackle','gk_diving','gk_kicking',
            'gk_positioning']

#, ,'sliding_tackle',  'gk_handling',,'gk_reflexes',,

target = 'overall_rating'

X = df[features]
y = df[target]

print(X.shape)
print(y.shape)

y = y.values.reshape(-1,1)


X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.3,random_state=2)

print(" The Training and Test size after splitting \n",'-'*80)
# size of the Training and Test set
print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)
```

```
In [ ]:   1  # Step 6  : Model /instantiating the Linear Regression Classifier
          2  lm = LinearRegression()
          3
          4  # fitting / Training data with the training set data
          5  lm.fit(X_train,y_train)
          6
          7  # Predicting the prices for the test set
          8  pred_test = lm.predict(X_test) # test case prediction
          9
         10  print(" Following are the predicted price values for the test set , after fitting the model with train set\n",'-'*80)
         11  print(pred_test)
         12
         13
```

```
In [ ]:   1  #Step 7 :  Visualize the predicted and observed ratings
          2
          3  # first, plot the observed data
          4
          5  plt.scatter(y_test, pred_test)
          6  plt.xlabel('Observed Rating')
          7  plt.ylabel('Predicted Overall Rating')
          8  plt.title('Actual vs. Predicted Overall Rating')
          9  plt.show()
         10
```

```
In [ ]:    1  # Step 8 : Model Evaluation
           2  #from sklearn.metrics import mean_squared_error,r2_score
           3
           4
           5
           6  RMSE = sqrt(mean_squared_error(y_test,pred_test))
           7  print("The Root Mean squared error for Model is : %.4f"%RMSE)
           8
           9  R2square = r2_score(y_test,pred_test)
          10
          11  print("The R2 score for the model is  : %.4f"%R2square)
          12
          13  acc = round(R2square,2)*100
          14
          15  print("\n This Model can predict the performance of the Player at %.2f percent accuracy" %acc)
          16
```