

```
In [1]: import pandas as pd

In [2]: import numpy as np

In [3]: import matplotlib.pyplot as plt

In [4]: import seaborn as sns

In [6]: df = pd.read_csv('https://github.com/YBI-Foundation/Dataset/raw/main/MPG.csv')

In [7]: df.head()

Out[7]:
   mpg  cylinders  displacement  horsepower  weight  acceleration  model_year  origin
0   18.0         8         307.0         130.0   3504          12.0         70     usa  chevrolet chevelle malibu
1   15.0         8         350.0         165.0   3693          11.5         70     usa    buick skylark 320
2   18.0         8         318.0         150.0   3436          11.0         70     usa    plymouth satellite
3   16.0         8         304.0         150.0   3433          12.0         70     usa      amc rebel sst
4   17.0         8         302.0         140.0   3449          10.5         70     usa      ford torino

In [8]: df.nunique()

Out[8]:
mpg           129
cylinders       5
displacement   82
horsepower     93
weight        351
acceleration   95
model_year     13
origin         3
name          385
dtype: int64

In [9]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 9 columns):
#   Column             Non-Null Count  Dtype
---  ---
0   mpg                 398 non-null   float64
1   cylinders           398 non-null   int64
2   displacement        398 non-null   float64
3   horsepower          392 non-null   float64
4   weight              398 non-null   int64
5   acceleration        398 non-null   float64
6   model_year          398 non-null   int64
7   origin              398 non-null   object
8   name                398 non-null   object
dtypes: float64(4), int64(3), object(2)
memory usage: 26.1+ KB

In [10]: df.describe()

Out[10]:
      mpg  cylinders  displacement  horsepower  weight  acceleration  model_year
count  398.000000   398.000000    398.000000   392.000000   398.000000   398.000000   398.000000
mean    23.514573    5.454774    193.425879   104.469388   2970.424623    15.568090    76.010050
std     7.815984    1.701004    104.269838    38.491160    846.841774    2.757689    3.697627
min     9.000000    3.000000    68.000000    46.000000   1613.000000    8.000000    70.000000
25%    17.500000    4.000000    104.250000    75.000000   2223.750000   13.825000    73.000000
50%    23.000000    4.000000    148.500000    93.500000   2803.500000   15.500000    76.000000
75%    29.000000    8.000000    262.000000   126.000000   3608.000000   17.175000    79.000000
max    46.600000    8.000000   455.000000   230.000000   5140.000000   24.800000    82.000000

In [11]: df.corr()

Out[11]:
      mpg  cylinders  displacement  horsepower  weight  acceleration  model_year
mpg      1.000000   -0.775396   -0.804203   -0.778427   -0.831741    0.420289    0.579267
cylinders -0.775396    1.000000    0.950721    0.842983    0.896017   -0.505419   -0.348746
displacement -0.804203  0.950721    1.000000    0.897257    0.932824   -0.543684   -0.370164
horsepower -0.778427  0.842983    0.897257    1.000000    0.864538   -0.689196   -0.416361
weight     -0.831741  0.896017    0.932824    0.864538    1.000000   -0.417457   -0.306564
acceleration 0.420289  -0.505419   -0.543684   -0.689196   -0.417457    1.000000    0.288137
model_year  0.579267  -0.348746   -0.370164   -0.416361   -0.306564    0.288137    1.000000

In [12]: df = df.dropna()

In [13]: df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 392 entries, 0 to 397
Data columns (total 9 columns):
#   Column             Non-Null Count  Dtype
---  ---
0   mpg                 392 non-null   float64
1   cylinders           392 non-null   int64
2   displacement        392 non-null   float64
3   horsepower          392 non-null   float64
4   weight              392 non-null   int64
5   acceleration        392 non-null   float64
6   model_year          392 non-null   int64
7   origin              392 non-null   object
8   name                392 non-null   object
dtypes: float64(4), int64(3), object(2)
memory usage: 30.6+ KB

In [14]: sns.pairplot(df, x_vars= ['displacement', 'horsepower', 'weight', 'acceleration'])

Out[14]:
<seaborn.axisgrid.PairGrid at 0x235076f25b0>

In [15]: sns.regplot(x = 'displacement', y = 'mpg', data = df);

In [16]: df.columns

Out[16]:
Index(['mpg', 'cylinders', 'displacement', 'horsepower', 'weight',
       'acceleration', 'model_year', 'origin', 'name'],
      dtype='object')

In [17]: y=df['mpg']

In [18]: y.shape

Out[18]:
(392,)

In [19]: X=df[['displacement','horsepower','weight','acceleration']]

In [20]: X.shape

Out[20]:
(392, 4)

In [21]: from sklearn.preprocessing import StandardScaler

In [22]: ss=StandardScaler()

In [23]: X=ss.fit_transform(X)

In [24]: X

Out[24]:
array([[ 1.07728956,  0.66413273,  0.62054034, -1.285258  ],
       [ 1.40873189,  1.57459447,  0.94333403, -1.40672302],
       [ 1.1825422 ,  1.18439658,  0.54038176, -1.64818924],
       ...,
       [-0.56847897, -0.53247413, -0.80463202, -1.4304305 ],
       [-0.7120053 , -0.66254089, -0.43562716,  1.13008813],
       [-0.72157372, -0.58450951, -0.30364091,  1.40043312]])

In [25]: from sklearn.model_selection import train_test_split

In [26]: X_train,X_test,y_train,y_test = train_test_split(X,y,train_size=0.7,random_state=72529)

In [27]: X_train.shape,X_test.shape,y_train.shape,y_test.shape

Out[27]:
((274, 4), (118, 4), (274,), (118,))

In [28]: from sklearn.linear_model import LinearRegression

In [29]: lr=LinearRegression()

In [30]: lr.fit(X_train,y_train)

Out[30]:
LinearRegression
LinearRegression()

In [31]: lr.intercept_

Out[31]:
23.446651543249907

In [32]: lr.coef_

Out[32]:
array([ 0.52918308, -1.88921858, -5.20327157,  0.22990503])

In [33]: y_pred = lr.predict(X_test)

In [34]: y_pred

Out[34]:
array([26.51044022, 28.32340153, 25.91553363, 31.72105637, 23.89426242,
       19.0751043 ,  7.59040685, 14.20914363, 27.66643119, 17.06330445,
       32.47072877, 15.73638657, 32.57265311, 30.93009747, 28.85653701,
       22.80700546, 24.39291745, 24.96816923, 33.38526647, 21.12723959,
       12.23062562, 15.90150476, 18.43202935, 23.37844768, 26.1573544 ,
       15.83303771, 23.8609154 , 24.35119494, 30.53150837, 27.27510839,
       24.12721904, 28.41176311, 27.08380862, 26.86573204, 30.10674388,
       31.01110316, 30.02863831, 30.76582169, 11.63349773,  8.68120334,
       23.26659201, 22.9920654 , 17.47928237, 27.59929556, 16.95268137,
       22.3102763 , 11.17098742, 30.77994215, 25.54865477, 26.20690975,
       23.10117752, 30.58074635, 24.54752048, 14.98070961, 16.93785991,
       9.73509395, 30.78860903, 23.50584069, 23.07395144, 26.20259207,
       14.02317006, 10.56156683, 31.29509943, 31.53533727, 15.33864478,
       22.93824319, 17.31513442, 12.88910908, 26.38150219, 10.57065935,
       17.48522681, 24.41979833, 32.58515765, 21.94919866, 23.03809382,
       28.74644542, 29.06515256, 24.25386401, 11.74091763, 22.02892 ,
       22.58484217, 25.99919669, 26.78187944, 23.86099505, 18.87111309,
       17.61153245, 10.93697163, 15.761219 , 16.53440609, 28.07300441,
       30.73808677, 27.05840352, 21.46236824, 10.18167757, 25.20415485,
       30.5827039 , 28.33624042, 25.85545963, 17.71534631, 24.09475365,
       29.69648792, 24.81344506, 16.94078205, 23.2359223 , 21.65705479,
       32.44175481, 26.38857705,  9.3320364 , 23.81535685, 24.29070372,
       30.56170743, 25.86200803, 21.07247079, 24.63791177, 26.96200317,
       24.47894339, 14.46848156, 28.93883923])

In [35]: from sklearn.metrics import mean_absolute_error, mean_absolute_percentage_error,r2_score

In [36]: mean_absolute_error(y_test,y_pred)

Out[36]:
3.406826049582863

In [37]: mean_absolute_percentage_error(y_test,y_pred)

Out[37]:
0.1554748737101835

In [38]: r2_score(y_test,y_pred)

Out[38]:
0.6969574372048439

In [39]: from sklearn.preprocessing import PolynomialFeatures

In [40]: poly = PolynomialFeatures(degree=2, interaction_only=True,include_bias=False)

In [41]: X_train2 = poly.fit_transform(X_train)

In [42]: X_test2 = poly.fit_transform(X_test)

In [43]: lr.fit(X_train2,y_train)

Out[43]:
LinearRegression
LinearRegression()

In [44]: lr.intercept_

Out[44]:
21.087825279131224

In [45]: lr.coef_

Out[45]:
array([-1.94714306, -5.75607397, -1.32768733, -0.60247703,  2.35766008,
        0.51842164,  0.63206027, -0.31391788, -0.39114212, -0.0318515 ])

In [46]: y_pred_poly = lr.predict(X_test2)

In [47]: from sklearn.metrics import mean_absolute_error, mean_absolute_percentage_error,r2_score

In [48]: mean_absolute_error(y_test,y_pred_poly)

Out[48]:
3.0481190089599766

In [49]: mean_absolute_percentage_error(y_test,y_pred_poly)

Out[49]:
0.13407106129461957

In [50]: r2_score(y_test,y_pred_poly)

Out[50]:
0.7281536246702335

In [ ]:
```