

Multivariate Regression Using Gradient descent

```
In [1]: #importing libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import math
import seaborn as sns
```

Importing Dataset

```
In [2]: dataset=pd.read_csv("Life Expectancy Data.csv")
```

```
In [3]: dataset.head()
```

Out[3]:

	Country	Year	Status	Life expectancy	Adult Mortality	infant deaths	Alcohol	percentage expenditure	Hepatit
0	Afghanistan	2015	Developing	65.0	263.0	62	0.01	71.279624	65
1	Afghanistan	2014	Developing	59.9	271.0	64	0.01	73.523582	62
2	Afghanistan	2013	Developing	59.9	268.0	66	0.01	73.219243	64
3	Afghanistan	2012	Developing	59.5	272.0	69	0.01	78.184215	67
4	Afghanistan	2011	Developing	59.2	275.0	71	0.01	7.097109	68

5 rows × 10 columns

```
In [4]: #displaying the dataset
dataset
```

Out[4]:

	Country	Year	Status	Life expectancy	Adult Mortality	infant deaths	Alcohol	percentage expenditure	Hep
0	Afghanistan	2015	Developing	65.0	263.0	62	0.01	71.279624	
1	Afghanistan	2014	Developing	59.9	271.0	64	0.01	73.523582	
2	Afghanistan	2013	Developing	59.9	268.0	66	0.01	73.219243	
3	Afghanistan	2012	Developing	59.5	272.0	69	0.01	78.184215	
4	Afghanistan	2011	Developing	59.2	275.0	71	0.01	7.097109	
...
2933	Zimbabwe	2004	Developing	44.3	723.0	27	4.36	0.000000	
2934	Zimbabwe	2003	Developing	44.5	715.0	26	4.06	0.000000	
2935	Zimbabwe	2002	Developing	44.8	73.0	25	4.43	0.000000	
2936	Zimbabwe	2001	Developing	45.3	686.0	25	1.72	0.000000	
2937	Zimbabwe	2000	Developing	46.0	665.0	24	1.68	0.000000	

2938 rows × 22 columns

Checking for NULL values

```
In [5]: pd.isna(dataset).any()
```

```
Out[5]: Country          False
Year          False
Status        False
Life expectancy    True
Adult Mortality    True
infant deaths     False
Alcohol           True
percentage expenditure  False
Hepatitis B       True
Measles          False
BMI              True
under-five deaths  False
Polio            True
Total expenditure    True
Diphtheria        True
HIV/AIDS         False
GDP              True
Population        True
thinness 1-19 years  True
thinness 5-9 years  True
Income composition of resources  True
Schooling         True
dtype: bool
```

```
In [6]: #checking details of the dataset  
dataset.describe()
```

Out[6]:

	Year	Life expectancy	Adult Mortality	infant deaths	Alcohol	percentage expenditure	Hep
count	2938.000000	2928.000000	2928.000000	2938.000000	2744.000000	2938.000000	2385
mean	2007.518720	69.224932	164.796448	30.303948	4.602861	738.251295	80
std	4.613841	9.523867	124.292079	117.926501	4.052413	1987.914858	25
min	2000.000000	36.300000	1.000000	0.000000	0.010000	0.000000	1
25%	2004.000000	63.100000	74.000000	0.000000	0.877500	4.685343	77
50%	2008.000000	72.100000	144.000000	3.000000	3.755000	64.912906	92
75%	2012.000000	75.700000	228.000000	22.000000	7.702500	441.534144	97
max	2015.000000	89.000000	723.000000	1800.000000	17.870000	19479.911610	99

Seeing the correlation matrix to select the feature

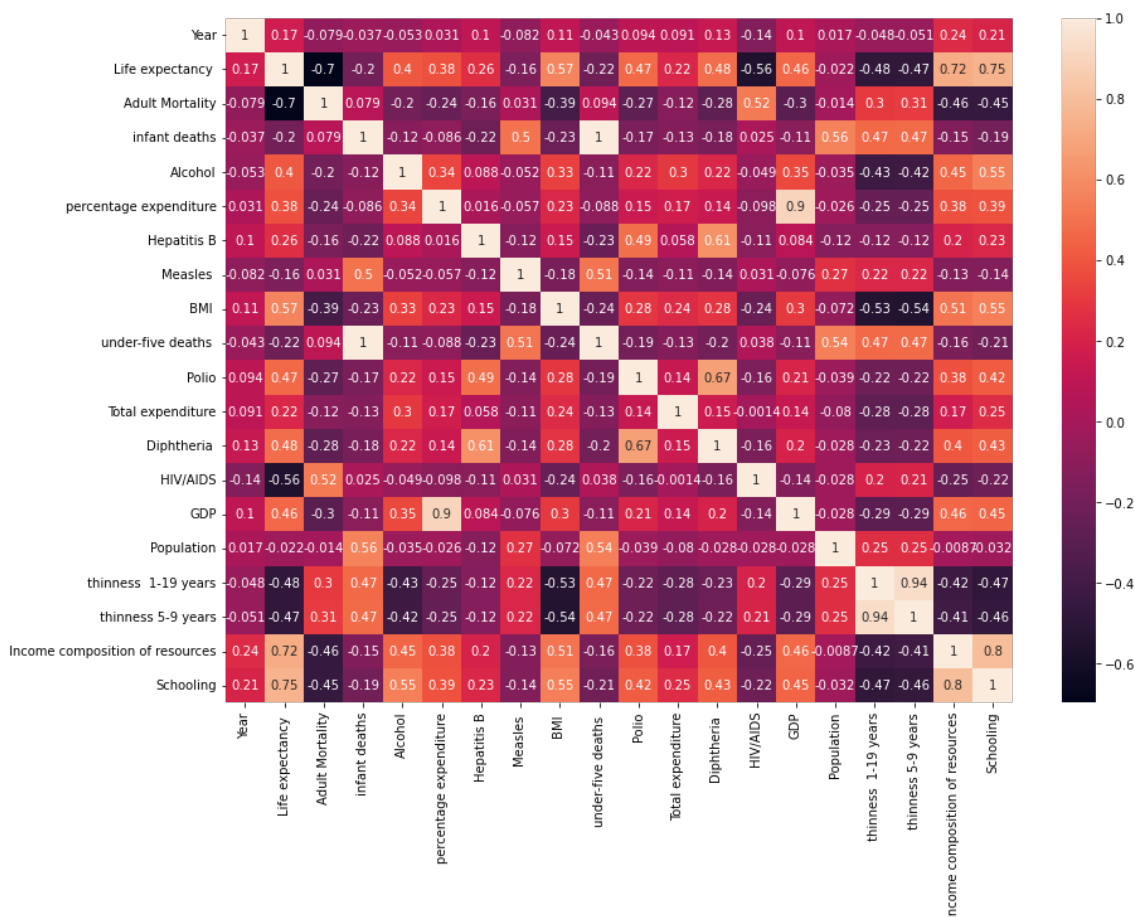
In [7]: `dataset.corr()`

Out[7]:

	Year	Life expectancy	Adult Mortality	infant deaths	Alcohol	percentage expenditure	Hepatitis B
Year	1.000000	0.170033	-0.079052	-0.037415	-0.052990	0.031400	0.104333
Life expectancy	0.170033	1.000000	-0.696359	-0.196557	0.404877	0.381864	0.256762
Adult Mortality	-0.079052	-0.696359	1.000000	0.078756	-0.195848	-0.242860	-0.162476
infant deaths	-0.037415	-0.196557	0.078756	1.000000	-0.115638	-0.085612	-0.223566
Alcohol	-0.052990	0.404877	-0.195848	-0.115638	1.000000	0.341285	0.087549
percentage expenditure	0.031400	0.381864	-0.242860	-0.085612	0.341285	1.000000	0.016274
Hepatitis B	0.104333	0.256762	-0.162476	-0.223566	0.087549	0.016274	1.000000
Measles	-0.082493	-0.157586	0.031176	0.501128	-0.051827	-0.056596	-0.120529
BMI	0.108974	0.567694	-0.387017	-0.227279	0.330408	0.228700	0.150380
under-five deaths	-0.042937	-0.222529	0.094146	0.996629	-0.112370	-0.087852	-0.233126
Polio	0.094158	0.465556	-0.274823	-0.170689	0.221734	0.147259	0.486171
Total expenditure	0.090740	0.218086	-0.115281	-0.128616	0.296942	0.174420	0.058280
Diphtheria	0.134337	0.479495	-0.275131	-0.175171	0.222020	0.143624	0.611495
HIV/AIDS	-0.139741	-0.556556	0.523821	0.025231	-0.048845	-0.097857	-0.112675
GDP	0.101620	0.461455	-0.296049	-0.108427	0.354712	0.899373	0.083903
Population	0.016969	-0.021538	-0.013647	0.556801	-0.035252	-0.025662	-0.123321
thinness 1-19 years	-0.047876	-0.477183	0.302904	0.465711	-0.428795	-0.251369	-0.120429
thinness 5-9 years	-0.050929	-0.471584	0.308457	0.471350	-0.417414	-0.252905	-0.124960
Income composition of resources	0.243468	0.724776	-0.457626	-0.145139	0.450040	0.381952	0.199549
Schooling	0.209400	0.751975	-0.454612	-0.193720	0.547378	0.389687	0.231117

```
In [8]: #Visuailizing the coorelation matrix
plt.figure(figsize = (14, 10))
correlation = dataset.corr()
sns.heatmap(correlation,
            xticklabels=correlation.columns,
            yticklabels=correlation.columns,
            annot=True)
```

Out[8]: <AxesSubplot:>



Defining Features and Target variable

```
In [9]: x=dataset[["Adult Mortality","infant deaths","Alcohol","percentage ex
y=dataset.iloc[:,3:4].values
```

In [10]: x

```
Out[10]: array([[2.63000000e+02, 6.20000000e+01, 1.00000000e-02, ...,
                5.84259210e+02, 4.79000000e-01, 1.01000000e+01],
                [2.71000000e+02, 6.40000000e+01, 1.00000000e-02, ...,
                6.12696514e+02, 4.76000000e-01, 1.00000000e+01],
                [2.68000000e+02, 6.60000000e+01, 1.00000000e-02, ...,
                6.31744976e+02, 4.70000000e-01, 9.90000000e+00],
                ...,
                [7.30000000e+01, 2.50000000e+01, 4.43000000e+00, ...,
                5.73483400e+01, 4.27000000e-01, 1.00000000e+01],
                [6.86000000e+02, 2.50000000e+01, 1.72000000e+00, ...,
                5.48587312e+02, 4.27000000e-01, 9.80000000e+00],
                [6.65000000e+02, 2.40000000e+01, 1.68000000e+00, ...,
                5.47358878e+02, 4.34000000e-01, 9.80000000e+00]])
```

In [11]: y

```
Out[11]: array([[65. ],
                [59.9],
                [59.9],
                ...,
                [44.8],
                [45.3],
                [46. ]])
```

Taking care of missing values

Using SimpleImputer Class of sklearn.impute library to fill the missing values

```
In [12]: # Using mean strategy to impute missing values
from sklearn.impute import SimpleImputer
imputer=SimpleImputer(missing_values=np.nan, strategy='mean')
imputer.fit(y[:, :])
y[:, :] = imputer.transform(y[:, :])
```

```
In [13]: from sklearn.impute import SimpleImputer
imputer=SimpleImputer(missing_values=np.nan, strategy='mean')
imputer.fit(x[:, :])
x[:, :] = imputer.transform(x[:, :])
```

In [14]: x

```
Out[14]: array([[2.63000000e+02, 6.20000000e+01, 1.00000000e-02, ...,
                5.84259210e+02, 4.79000000e-01, 1.01000000e+01],
                [2.71000000e+02, 6.40000000e+01, 1.00000000e-02, ...,
                6.12696514e+02, 4.76000000e-01, 1.00000000e+01],
                [2.68000000e+02, 6.60000000e+01, 1.00000000e-02, ...,
                6.31744976e+02, 4.70000000e-01, 9.90000000e+00],
                ...,
                [7.30000000e+01, 2.50000000e+01, 4.43000000e+00, ...,
                5.73483400e+01, 4.27000000e-01, 1.00000000e+01],
                [6.86000000e+02, 2.50000000e+01, 1.72000000e+00, ...,
                5.48587312e+02, 4.27000000e-01, 9.80000000e+00],
                [6.65000000e+02, 2.40000000e+01, 1.68000000e+00, ...,
                5.47358878e+02, 4.34000000e-01, 9.80000000e+00]])
```

In [15]:

y

```
Out[15]: array([[65. ],
               [59.9],
               [59.9],
               ...,
               [44.8],
               [45.3],
               [46. ]])
```

In [16]:

len(x)

Out[16]: 2938

In [17]:

len(y)

Out[17]: 2938

In [18]:

x.shape, y.shape

Out[18]: ((2938, 13), (2938, 1))

Feature Scaling

```
In [19]: from sklearn.preprocessing import StandardScaler
sc_x=StandardScaler()
sc_y=StandardScaler()
x=sc_x.fit_transform(x)
y=sc_y.fit_transform(y)
```

$$X_{\text{new}} = \frac{X_i - X_{\text{mean}}}{\text{Standard Deviation}}$$

```
In [20]: #setting the matrixes
#adding extra column(of value 1) in the feature for the symmetry. Do
#the matrix multiplication between our transpose of parameter vector
x = x[:,0:len(x)]
ones = np.ones([x.shape[0],1])
x = np.concatenate((ones,x),axis=1)

#y = y.iloc[:,:].values#.values converts it from pandas.core.frame.Da
theta = np.zeros([1,14])
```

In [21]:

x.shape

Out[21]: (2938, 14)

Splitting Dataset into Training Set and Test Set

```
In [22]: x_test=x[0:587,:]  
x_train=x[587,:]  
y_test=y[0:587,:]  
y_train=y[587,:]
```

```
In [23]: x_train
```

```
Out[23]: array([[ 1.          , -0.13539061, -0.1213161 , ..., -0.54885025,  
                  0.14868742, -0.21226403],  
                [ 1.          , -1.20745957, -0.1213161 , ..., -0.39871143,  
                  0.15357059, -0.12034736],  
                [ 1.          , -0.01448057, -0.1128348 , ..., -0.39037705,  
                  0.13892107, -0.15098625],  
                ...,  
                [ 1.          , -0.73994077, -0.04498439, ..., -0.56536401,  
                  -0.97932554, -0.61056961],  
                [ 1.          ,  4.20124926, -0.04498439, ..., -0.52796354,  
                  -0.97932554, -0.6718474 ],  
                [ 1.          ,  4.03197521, -0.05346569, ..., -0.52805706,  
                  -0.94514333, -0.6718474 ]])
```

```
In [24]: x_test
```

```
Out[24]: array([[ 1.          ,  0.79158632,  0.26882378, ..., -0.52524766,  
                  -0.72540055, -0.57993072],  
                [ 1.          ,  0.85607167,  0.28578638, ..., -0.52308258,  
                  -0.74005007, -0.61056961],  
                [ 1.          ,  0.83188966,  0.30274898, ..., -0.52163233,  
                  -0.7693491 , -0.6412085 ],  
                ...,  
                [ 1.          , -0.16763328, -0.1382787 , ..., -0.21385854,  
                  0.23170136, -0.02843068],  
                [ 1.          , -0.14345128, -0.1297974 , ..., -0.54081509,  
                  0.20240232, -0.15098625],  
                [ 1.          , -0.16763328, -0.1297974 , ..., -0.31191765,  
                  0.14868742, -0.27354181]])
```

```
In [25]: y_train
```

```
Out[25]: array([[ 0.37608459],  
                [ 0.334006  ],  
                [ 0.27088811],  
                ...,  
                [-2.56941673],  
                [-2.5168185  ],  
                [-2.44318096]])
```



```
Out[26]: array([[ -0.44444792],  
                [ -0.98094995],  
                [ -0.98094995],  
                [-1.02302854],  
                [-1.05458748],  
                [-1.09666607],  
                [-1.11770537],  
                [-1.17030361],  
                [-1.23342149],  
                [-1.25446079],  
                [-1.25446079],  
                [-1.28601973],  
                [-1.31757867],  
                [-1.37017691],  
                [-1.46485374],  
                [-1.51745198],  
                [ 0.90206696],  
                [ 0.87050802],  
                [ 0.83894908],  
                [ 0.80730014]])
```

```
In [27]: #compute cost
def Cost(X,y,theta):
    tobesummed = np.power((X @ theta.T - y),2)
    return np.nansum(tobesummed)/(2 * len(X))
def gradientDescent(X,y,theta,iters,alpha):
    cost = np.zeros(iters)
    for i in range(iters):
        theta = theta - (alpha/len(X)) * np.nansum(X * (X @ theta.T - y))
        #print(theta)
        cost[i] = Cost(X, y, theta)

    return theta,cost
```

$$\frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

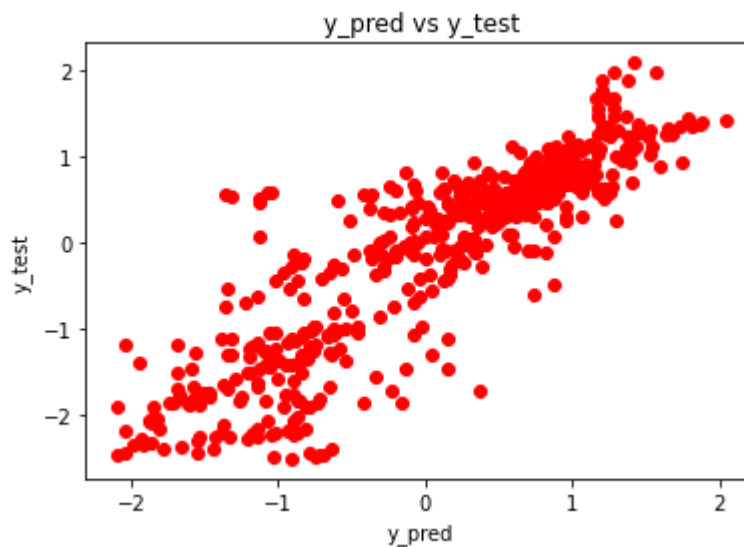
$$\begin{aligned} & \text{repeat until convergence } \{ \\ & \quad \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \\ & \quad \text{(for } j = 1 \text{ and } j = 0) \\ & \} \end{aligned}$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

```
[[ 0.0114838 -0.38513849 -0.00392374  0.02422584  0.00290746 -0.02
712843
  -0.06872925  0.11919904  0.0645791  0.01874079  0.10268727  0.05
955247
   0.19021389  0.22495502]]
0.11758367259575796
```

```
Out[31]: array([[ -8.63872084e-01],  
                [ -7.62789897e-01],  
                [ -7.51582507e-01],  
                [ -7.68129568e-01],  
                [ -8.12314370e-01],  
                [ -8.47040513e-01],  
                [ -9.11778927e-01],  
                [ -9.43498464e-01],  
                [ -1.02454645e+00],  
                [ -1.09431635e+00],  
                [ -1.09659827e+00],  
                [ -1.56649249e+00],  
                [ -1.34674060e+00],  
                [ -5.45848980e-01],  
                [ -1.59378402e+00],  
                [ -1.69297318e+00],  
                [  9.01368913e-01],  
                [  1.09605753e+00],  
                [  8.59082853e-01],  
                [  9.42140074e-01]]
```

```
In [32]: #Plotting the scatter plot to analyze test target variable and prediction
plt.scatter(y_pred,y_test,color='red')
plt.title('y_pred vs y_test')
plt.xlabel('y_pred')
plt.ylabel("y_test")
plt.show()
```



```
In [33]: print(np.concatenate((y_pred.reshape(1,len(y_pred)),y_test.reshape(1,
[[ -0.86387208 -0.7627899 -0.75158251 ...  0.37910093  0.31277057
    0.01865373]
 [ -0.44444792 -0.98094995 -0.98094995 ...  0.44972212  0.40764353
    0.40764353])
```

Defining the RMSE

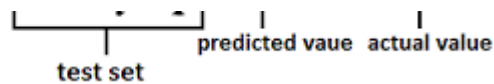
$$RMSE = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$

```
In [34]: def rmse(predictions, targets):
          return np.sqrt(((predictions - targets) ** 2).mean())
rmse(y_pred, y_test)
```

Out[34]: 0.5470236141843329

Defining the MAE

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$



```
In [35]: def MAE(y_pred,y):
          res=np.sum(abs(y_pred-y))/len(y)
          return res
          MAE(y_pred, y_test)
```

Out[35]: 0.39549713195690395

Multivariate Regression by Closed Form

Defining feature and target variable

```
In [36]: x1=dataset[["Adult Mortality","infant deaths","Alcohol","percentage of population with access to electricity"]]
          y1=dataset.iloc[:,3:4].values
```

```
In [37]: x
```

```
Out[37]: array([[ 1.          ,  0.79158632,  0.26882378, ..., -0.52524766,
                 -0.72540055, -0.57993072],
                [ 1.          ,  0.85607167,  0.28578638, ..., -0.52308258,
                 -0.74005007, -0.61056961],
                [ 1.          ,  0.83188966,  0.30274898, ..., -0.52163233,
                 -0.7693491 , -0.6412085 ],
                ...,
                [ 1.          , -0.73994077, -0.04498439, ..., -0.56536401,
                 -0.97932554, -0.61056961],
                [ 1.          ,  4.20124926, -0.04498439, ..., -0.52796354,
                 -0.97932554, -0.6718474 ],
                [ 1.          ,  4.03197521, -0.05346569, ..., -0.52805706,
                 -0.94514333, -0.6718474 ]])
```

```
In [38]: y
```

```
Out[38]: array([[ -0.44444792],
                [ -0.98094995],
                [ -0.98094995],
                ...,
                [-2.56941673],
                [-2.5168185 ],
                [-2.44318096]])
```

Taking care of missing values

```
In [39]: from sklearn.impute import SimpleImputer
          imputer=SimpleImputer(missing_values=np.nan,strategy='mean')
          imputer.fit(y[:,:])
          y[:,:]=imputer.transform(y[:,:])
```

```
In [40]: from sklearn.impute import SimpleImputer
imputer=SimpleImputer(missing_values=np.nan, strategy='mean')
imputer.fit(x[:, :])
x[:, :]=imputer.transform(x[:, :])
```

Feature Scaling

```
In [41]: from sklearn.preprocessing import StandardScaler
sc_x=StandardScaler()
sc_y=StandardScaler()
x1=sc_x.fit_transform(x1)
y1=sc_y.fit_transform(y1)
```

$$X_{\text{new}} = \frac{X_i - X_{\text{mean}}}{\text{Standard Deviation}}$$

```
In [42]: x1.shape, theta.shape, y1.shape
```

```
Out[42]: ((2938, 13), (1, 14), (2938, 1))
```

Splitting Dataset into Training set and Test set

```
In [43]: x1_test=x[0:587,:]  
x1_train=x[587:,:]  
y1_test=y[0:587,:]  
y1_train=y[587:,:]
```

Training model on training set

```
In [44]: def find_theta(x, y):

    m = x.shape[0] # Number of training examples.
    # Appending a cloumn of ones in X to add the bias term.
    x = np.append(x, np.ones((m,1)), axis=1)
    # reshaping y to (m,1)
    y = y.reshape(m,1)

    # The Normal Equation
    theta = np.dot(np.linalg.inv(np.dot(x.T, x)), np.dot(x.T, y))

    return theta
```

$$\theta = (X^T X)^{-1} X^T y$$

$$W = (X^T X)^{-1} X^T y$$

```
In [45]: def predict(x, theta):

    # Appending a cloumn of ones in X to add the bias term.
    x = np.append(x, np.ones((x.shape[0],1)), axis=1)

    # preds is y_hat which is the dot product of X and theta.
    preds = np.dot(x, theta)

    return preds
```

Predicting value on test set

```
In [46]: theta1 = find_theta(x1_train, y1_train)
         theta1
         y1_pred = predict(x1_test, theta1)
```

```
In [47]: y1_pred
```

```
Out[47]: array([[ -0.75],  
                [ -0.75],  
                [ -0.75],  
                [ -0.75],  
                [ -0.75],  
                [-1.   ],  
                [-1.   ],  
                [-1.   ],  
                [-1.   ],  
                [-1.   ],  
                [-1.   ],  
                [-1.75],  
                [-1.75],  
                [-1.   ],  
                [-2.   ],  
                [-2.   ],  
                [ 0.5  ],  
                [ 0.75],  
                [ 0.5  ],  
                [ 0.5  ]])
```

Defining RMSE

```
In [48]: def rmse(predictions, targets):
          return np.sqrt(((predictions - targets) ** 2).mean())
          rmse(y1_pred, y1_test)
```

Out[48]: 0.5819698495897416

Defining MAE

```
In [49]: def MAE(y_pred,y):  
         res=np.sum(abs(y_pred-y))/len(y)  
         return res  
MAE(y1_pred, y1_test)
```

Out[49]: 0.4362348846476966

```
In [50]: #plotting to analyze the y_pred and y_actual values  
plt.scatter(y1_pred,y1_test,color='red') #plotting the actual data  
plt.title('y_pred vs y_test')  
plt.xlabel('y_pred')  
plt.ylabel('y_test')  
plt.show()
```

