

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
from google.colab import files
upload=files.upload()
df = pd.read_csv('abalone.csv')
```

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving abalone.csv to abalone.csv

```
df.describe()
```

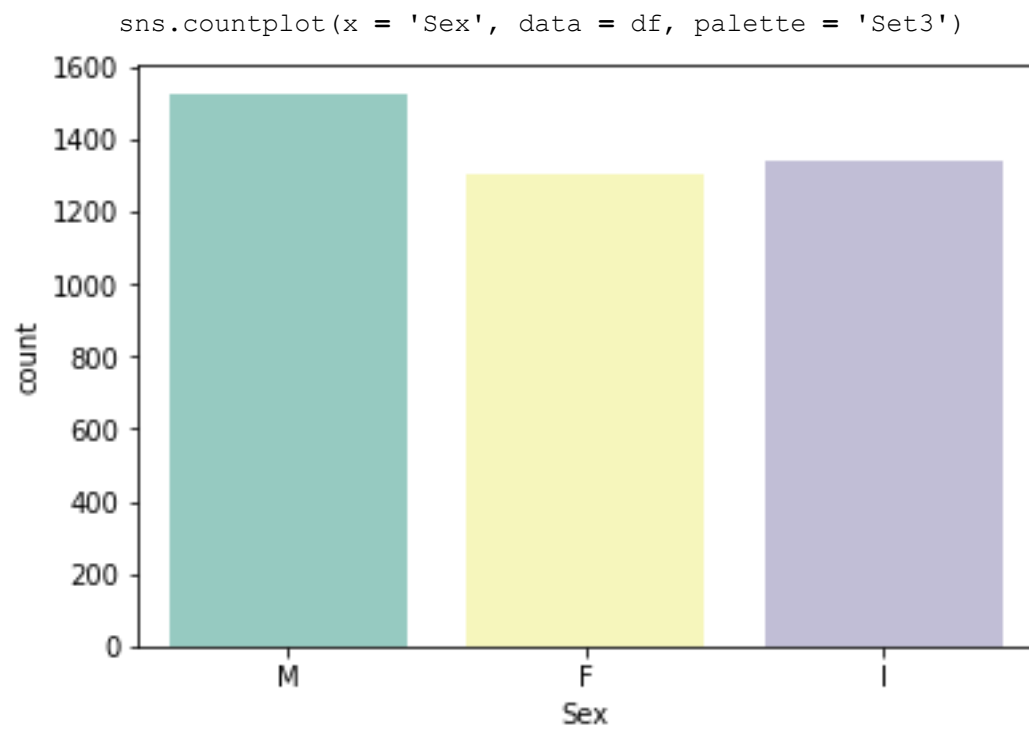
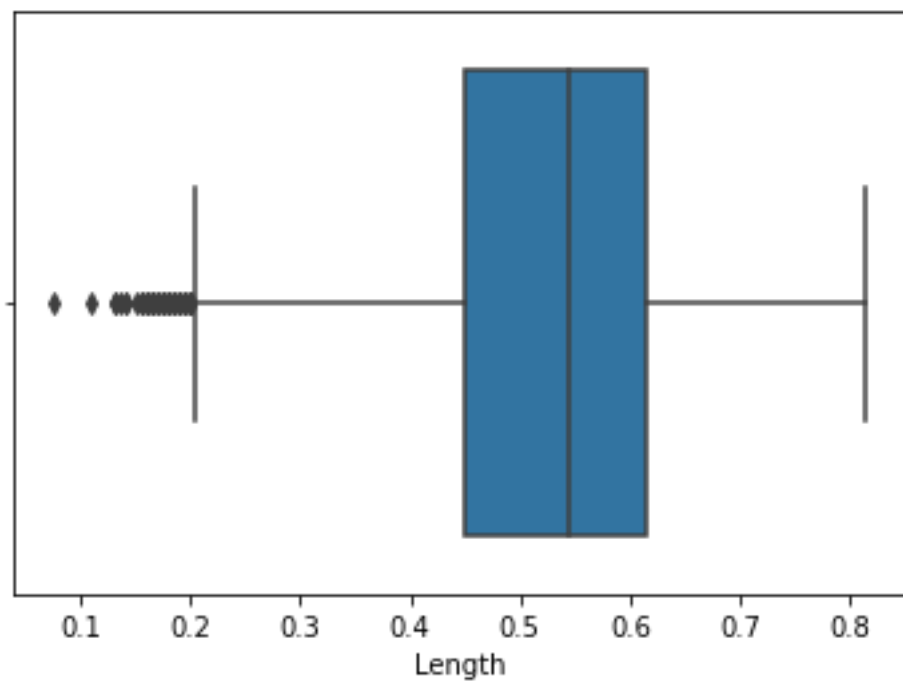
	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	Rings
count	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000
mean	0.523992	0.407881	0.139516	0.828742	0.359367	0.180594	0.238831	9.933684
std	0.120093	0.099240	0.041827	0.490389	0.221963	0.109614	0.139203	3.224169
min	0.075000	0.055000	0.000000	0.002000	0.001000	0.000500	0.001500	1.000000
25%	0.450000	0.350000	0.115000	0.441500	0.186000	0.093500	0.130000	8.000000
50%	0.545000	0.425000	0.140000	0.799500	0.336000	0.171000	0.234000	9.000000
75%	0.615000	0.480000	0.165000	1.153000	0.502000	0.253000	0.329000	11.000000
max	0.815000	0.650000	1.130000	2.825500	1.488000	0.760000	1.005000	29.000000

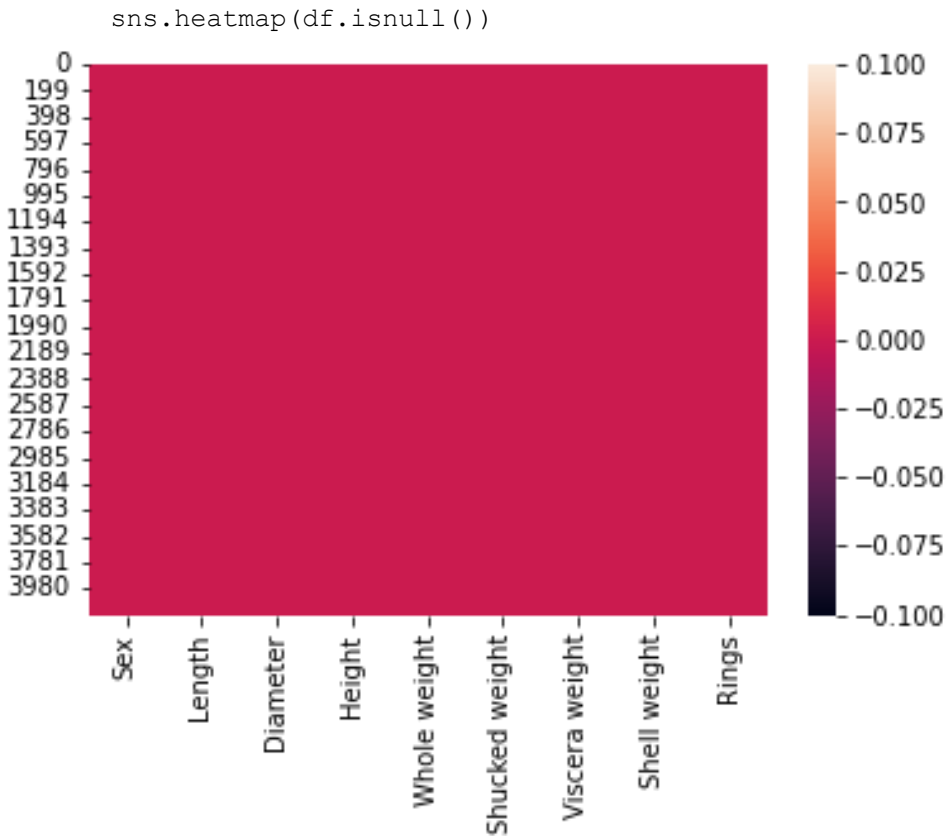
```
df.head()
```

	Sex	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	Rings
0	M	0.455	0.365	0.095	0.5140	0.2245	0.1010	0.150	15
1	M	0.350	0.265	0.090	0.2255	0.0995	0.0485	0.070	7
2	F	0.530	0.420	0.135	0.6770	0.2565	0.1415	0.210	9
3	M	0.440	0.365	0.125	0.5160	0.2155	0.1140	0.155	10
4	I	0.330	0.255	0.080	0.2050	0.0895	0.0395	0.055	7

```
#Perform visualisations
#Univariate analysis
```

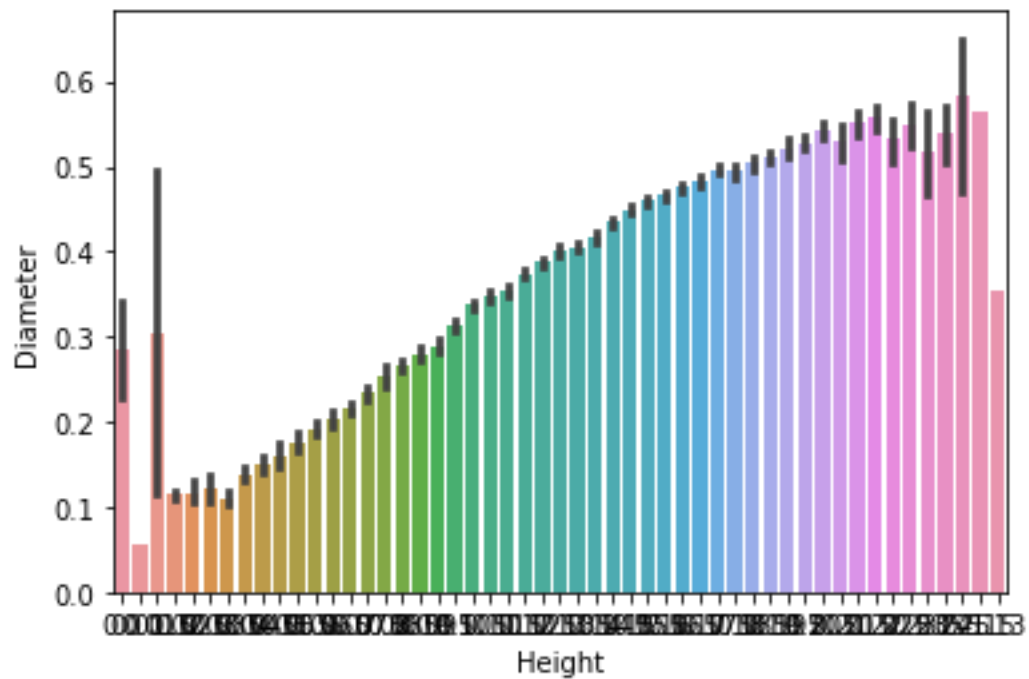
```
sns.boxplot(df.Length)
```





#Bivariate analysis

```
sns.barplot(x=df.Height,y=df.Diameter)
```

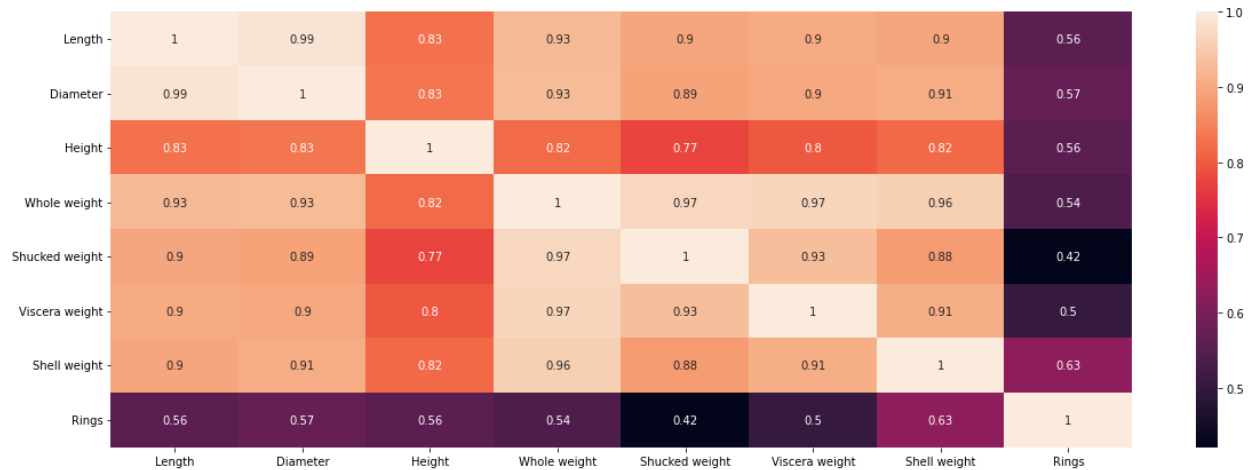


```
numerical_features = df.select_dtypes(include = [np.number]).columns
```

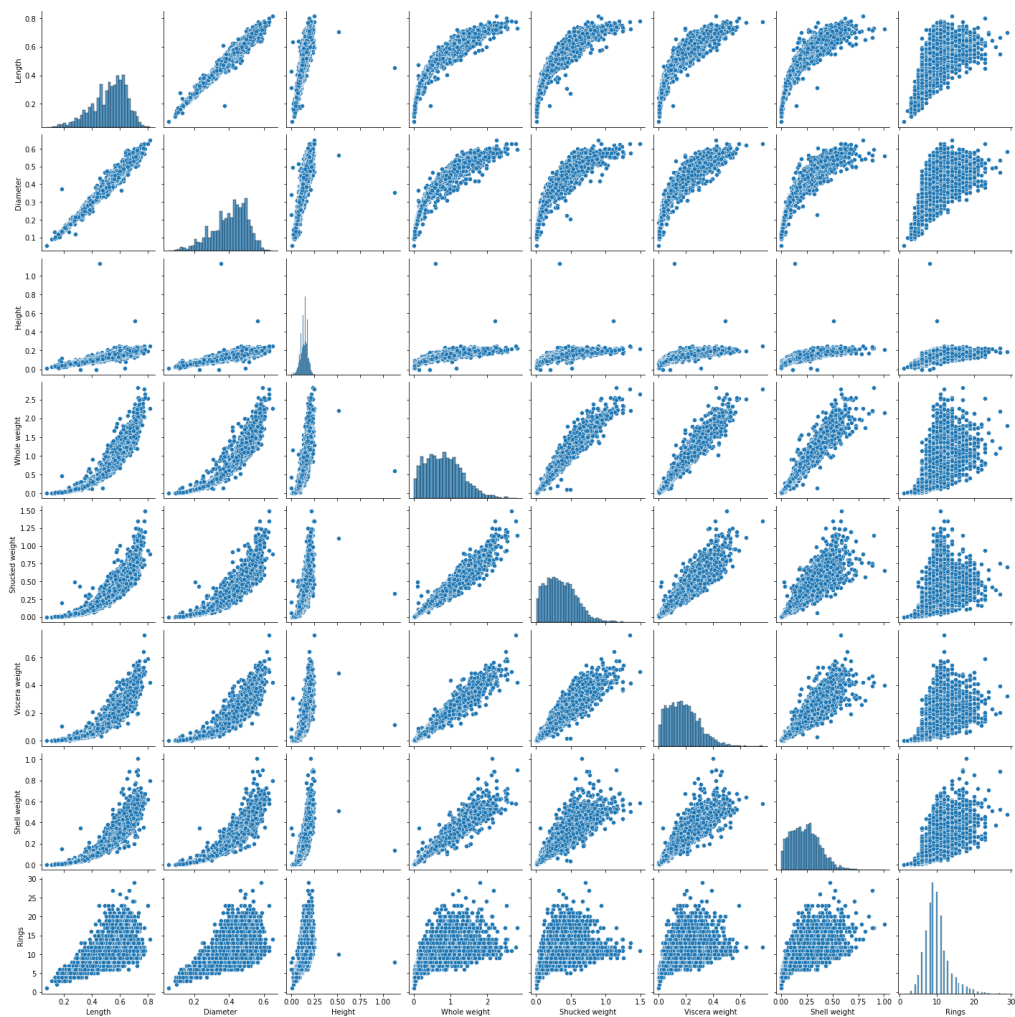
```
categorical_features = df.select_dtypes(include = [np.object]).columns
```

```
plt.figure(figsize = (20,7))
```

```
sns.heatmap(df[numerical_features].corr(),annot = True)
```



```
#Multivariate Analysis  
sns.pairplot(df)
```



```
#Perform descriptive model on the dataset
df['Height'].describe()
```

```
count      4177.000000
mean        0.139516
std         0.041827
min         0.000000
25%         0.115000
50%         0.140000
75%         0.165000
max         1.130000
Name: Height, dtype: float64
```

```
df['Height'].mean()
```

```
0.13951639932966242
```

```
df.max()
```

```
Sex          M
Length       0.815
Diameter     0.65
Height       1.13
Whole weight 2.8255
Shucked weight 1.488
Viscera weight 0.76
Shell weight 1.005
Rings        29
dtype: object
```

```
df['Sex'].value_counts()
```

```
M      1528
I      1342
F      1307
Name: Sex, dtype: int64
```

```
df[df.Height == 0]
```

	Sex	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	Rings
1257	I	0.430	0.34	0.0	0.428	0.2065	0.0860	0.1150	8
3996	I	0.315	0.23	0.0	0.134	0.0575	0.0285	0.3505	6

```
df['Shucked weight'].kurtosis()
```

```
0.5951236783694207
```

```
df['Diameter'].median()
```

```
0.425
```

```
df['Shucked weight'].skew()
```

```
0.7190979217612694
```

```
#Missing values
```

```
df.isna().any()
```

```
Sex                False
Length            False
Diameter           False
Height            False
Whole weight       False
Shucked weight     False
Viscera weight     False
Shell weight       False
Rings              False
dtype: bool
```

	Missing values	% Missing
Sex	0	0.0
Length	0	0.0
Diameter	0	0.0
Height	0	0.0
Whole weight	0	0.0
Shucked weight	0	0.0
Viscera weight	0	0.0
Shell weight	0	0.0
Rings	0	0.0

```

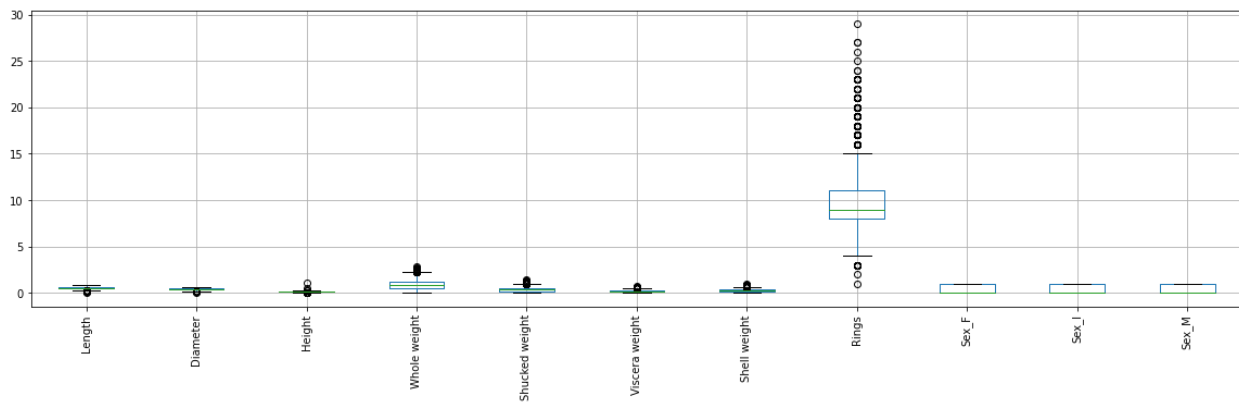
#Find the outliers
q1=df.Rings.quantile(0.25)
q2=df.Rings.quantile(0.75)
iqr=q2-q1
print(iqr)
3.0

```

```

df = pd.get_dummies(df)
dummy_df = df
df.boxplot( rot = 90, figsize=(20,5))

```



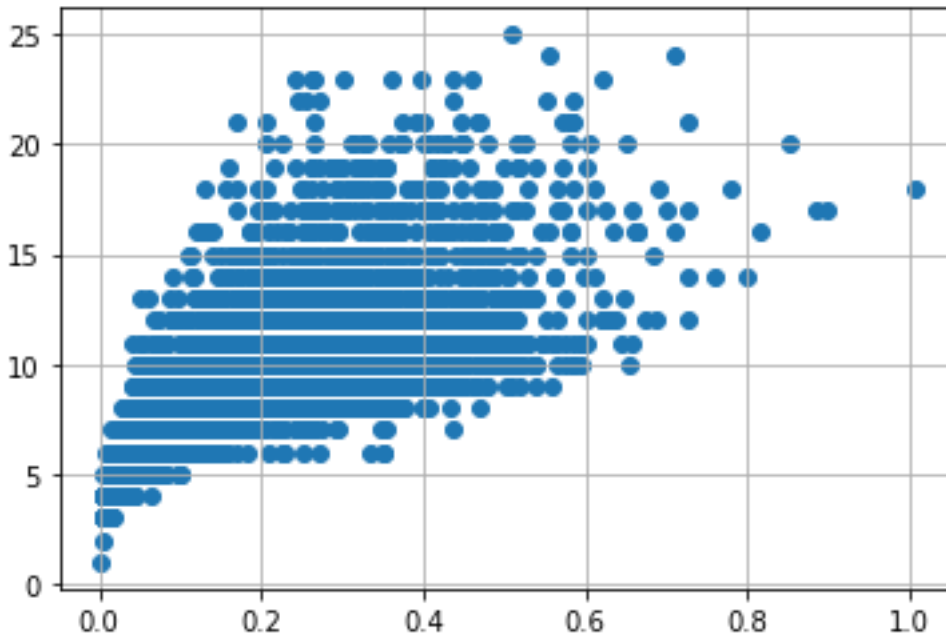
```

df['age'] = df['Rings']
df = df.drop('Rings', axis = 1)

df.drop(df[(df['Viscera weight']> 0.5) & (df['age'] < 20)].index,
inplace=True)
df.drop(df[(df['Viscera weight']<0.5) & (df['age'] > 25)].index,
inplace=True)

var = 'Shell weight'
plt.scatter(x = df[var], y = df['age'])
plt.grid(True)

```



#Check for categorical columns and perform encoding

```
numerical_features = df.select_dtypes(include = [np.number]).columns
categorical_features = df.select_dtypes(include = [np.object]).columns
```

```
numerical_features
categorical_features
```

```
Index([], dtype='object')
```

```
abalone_numeric = df[['Length', 'Diameter', 'Height', 'Whole weight',
'Shucked weight', 'Viscera weight', 'Shell weight', 'age', 'Sex_F', 'Sex_I',
'Sex_M']]
```

```
abalone_numeric.head()
```

	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	age	Sex_F	Sex_I	Sex_M
0	0.455	0.365	0.095	0.5140	0.2245	0.1010	0.150	15	0	0	1
1	0.350	0.265	0.090	0.2255	0.0995	0.0485	0.070	7	0	0	1
2	0.530	0.420	0.135	0.6770	0.2565	0.1415	0.210	9	1	0	0
3	0.440	0.365	0.125	0.5160	0.2155	0.1140	0.155	10	0	0	1
4	0.330	0.255	0.080	0.2050	0.0895	0.0395	0.055	7	0	1	0

#Dependent and Independent Variables

```
x = df.iloc[:, 0:1].values
```

```
y = df.iloc[:, 1]
```



```

y
0      0.365
1      0.265
2      0.420
3      0.365
4      0.255
...
4172    0.450
4173    0.440
4174    0.475
4175    0.485
4176    0.555
Name: Diameter, Length: 4150, dtype: float64

```

#Scaling the Independent Variables

```
print ("\n ORIGINAL VALUES: \n\n", x,y)
```

ORIGINAL VALUES:

```

[[0.455]
 [0.35 ]
 [0.53 ]
 ...
 [0.6  ]
 [0.625]
 [0.71 ]] 0      0.365
1      0.265
2      0.420
3      0.365
4      0.255
...
4172    0.450
4173    0.440
4174    0.475
4175    0.485
4176    0.555

```

Name: Diameter, Length: 4150, dtype: float64

```
from sklearn import preprocessing
```

```
min_max_scaler = preprocessing.MinMaxScaler(feature_range =(0, 1))
```

```
new_y= min_max_scaler.fit_transform(x,y)
```

```
print ("\n VALUES AFTER MIN MAX SCALING: \n\n", new_y)
```

VALUES AFTER MIN MAX SCALING:

```

[[0.51351351]
 [0.37162162]
 [0.61486486]
 ...
 [0.70945946]
 [0.74324324]
 [0.85810811]]

```

#Split the data into Training and Testing

```
X = df.drop('age', axis = 1)
```

```
y = df['age']
```

```

from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.feature_selection import SelectKBest
standardScale = StandardScaler()
standardScale.fit_transform(X)

selectkBest = SelectKBest()
X_new = selectkBest.fit_transform(X, y)

X_train, X_test, y_train, y_test = train_test_split(X_new, y, test_size =
0.25)
X_train
array([[0.525, 0.41 , 0.135, ..., 1.    , 0.    , 0.    ],
       [0.275, 0.175, 0.09 , ..., 0.    , 1.    , 0.    ],
       [0.68 , 0.56 , 0.195, ..., 1.    , 0.    , 0.    ],
       ...,
       [0.55 , 0.425, 0.15 , ..., 0.    , 1.    , 0.    ],
       [0.35 , 0.26 , 0.09 , ..., 0.    , 0.    , 1.    ],
       [0.57 , 0.42 , 0.14 , ..., 0.    , 0.    , 1.    ]])
y_train
2983      8
1764      5
888       11
2029      9
3096      9
...
279       11
584       11
581       14
3315      9
2835      8
Name: age, Length: 3112, dtype: int64
# Build the model
# Linear Regression

from sklearn import linear_model as lm
from sklearn.linear_model import LinearRegression
model=lm.LinearRegression()
results=model.fit(X_train,y_train)

accuracy = model.score(X_train, y_train)
print('Accuracy of the model:', accuracy)
Accuracy of the model: 0.5354279264706927
#Training the model
lm = LinearRegression()
lm.fit(X_train, y_train)
y_train_pred = lm.predict(X_train)
y_train_pred
array([ 8.203125,  6.34375 , 11.046875, ...,  9.359375,  8.09375 ,
        9.90625 ])
X_train
array([[0.525, 0.41 , 0.135, ..., 1.    , 0.    , 0.    ],

```

```

        [0.275, 0.175, 0.09 , ..., 0.    , 1.    , 0.    ],
        [0.68 , 0.56 , 0.195, ..., 1.    , 0.    , 0.    ],
        ...,
        [0.55 , 0.425, 0.15 , ..., 0.    , 1.    , 0.    ],
        [0.35 , 0.26 , 0.09 , ..., 0.    , 0.    , 1.    ],
        [0.57 , 0.42 , 0.14 , ..., 0.    , 0.    , 1.    ]])
y_train
2983      8
1764      5
888       11
2029      9
3096      9
..
279       11
584       11
581       14
3315      9
2835      8
Name: age, Length: 3112, dtype: int64
from sklearn.metrics import mean_absolute_error, mean_squared_error
s = mean_squared_error(y_train, y_train_pred)
print('Mean Squared error of training set :%2f'%s)

Mean Squared error of training set :4.696701
#Testing the model

y_train_pred = lm.predict(X_train)
y_test_pred = lm.predict(X_test)

y_test_pred
array([ 7.125    ,  9.        , 10.59375 , ...,  8.15625 ,  7.078125,
        9.609375])

X_test
array([[0.35 , 0.26 , 0.095, ..., 0.    , 1.    , 0.    ],
       [0.53 , 0.42 , 0.17 , ..., 1.    , 0.    , 0.    ],
       [0.525, 0.425, 0.16 , ..., 1.    , 0.    , 0.    ],
       ...,
       [0.35 , 0.265, 0.11 , ..., 0.    , 0.    , 1.    ],
       [0.425, 0.34 , 0.105, ..., 0.    , 1.    , 0.    ],
       [0.605, 0.47 , 0.165, ..., 0.    , 0.    , 1.    ]])

y_test
3813      8
2581      6
49        9
384       10
3832      14
..
3065      11
724       11
2311      7
1444      6
1006      11
Name: age, Length: 1038, dtype: int64

```

```
p = mean_squared_error(y_test, y_test_pred)
print('Mean Squared error of testing set :%2f'%p)
Mean Squared error of testing set :4.994425
```

#Measure the performance using metrics

```
from sklearn.metrics import r2_score
s = r2_score(y_train, y_train_pred)
print('R2 Score of training set:%.2f'%s)
```

R2 Score of training set:0.54

```
from sklearn.metrics import r2_score
p = r2_score(y_test, y_test_pred)
print('R2 Score of testing set:%.2f'%p)
```

R2 Score of testing set:0.51