

## IMPORTING LIBRARIES

```
In [1]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression
```

## 2. Load the dataset into the Google Colab

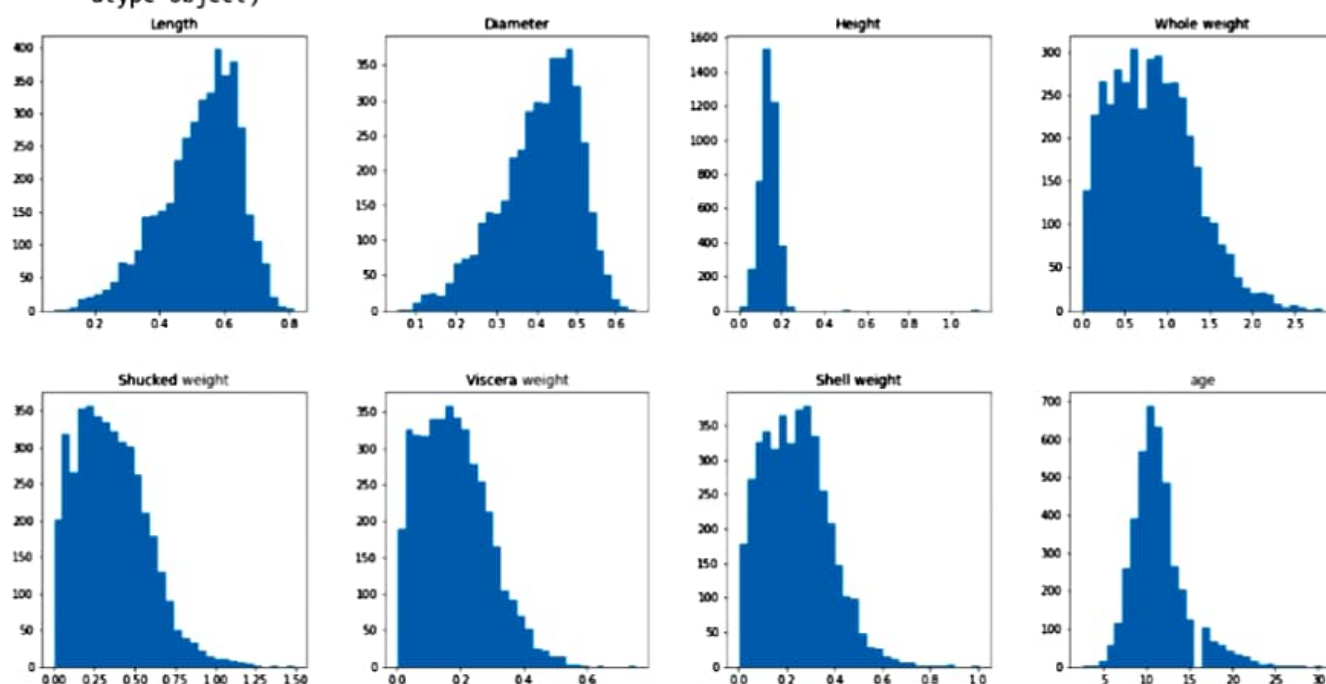
```
In [2]: df=pd.read_csv("/content/abalone.csv")
```

```
In [3]: df['age'] = df['Rings']+1.5
df = df.drop('Rings', axis = 1)
```

## 3. UNIVARIATE ANALYSIS

```
In [4]: df.hist(figsize=(20,10), grid=False, layout=(2, 4), bins = 30)
```

```
Out[4]: array([[
,
],
[,
,
],
],
dtype=object)
```



```
In [5]: df.groupby('Sex')[['Length', 'Diameter', 'Height', 'Whole weight', 'Shucked weight',
'Viscera weight', 'Shell weight', 'age']].mean().sort_values('age')
```

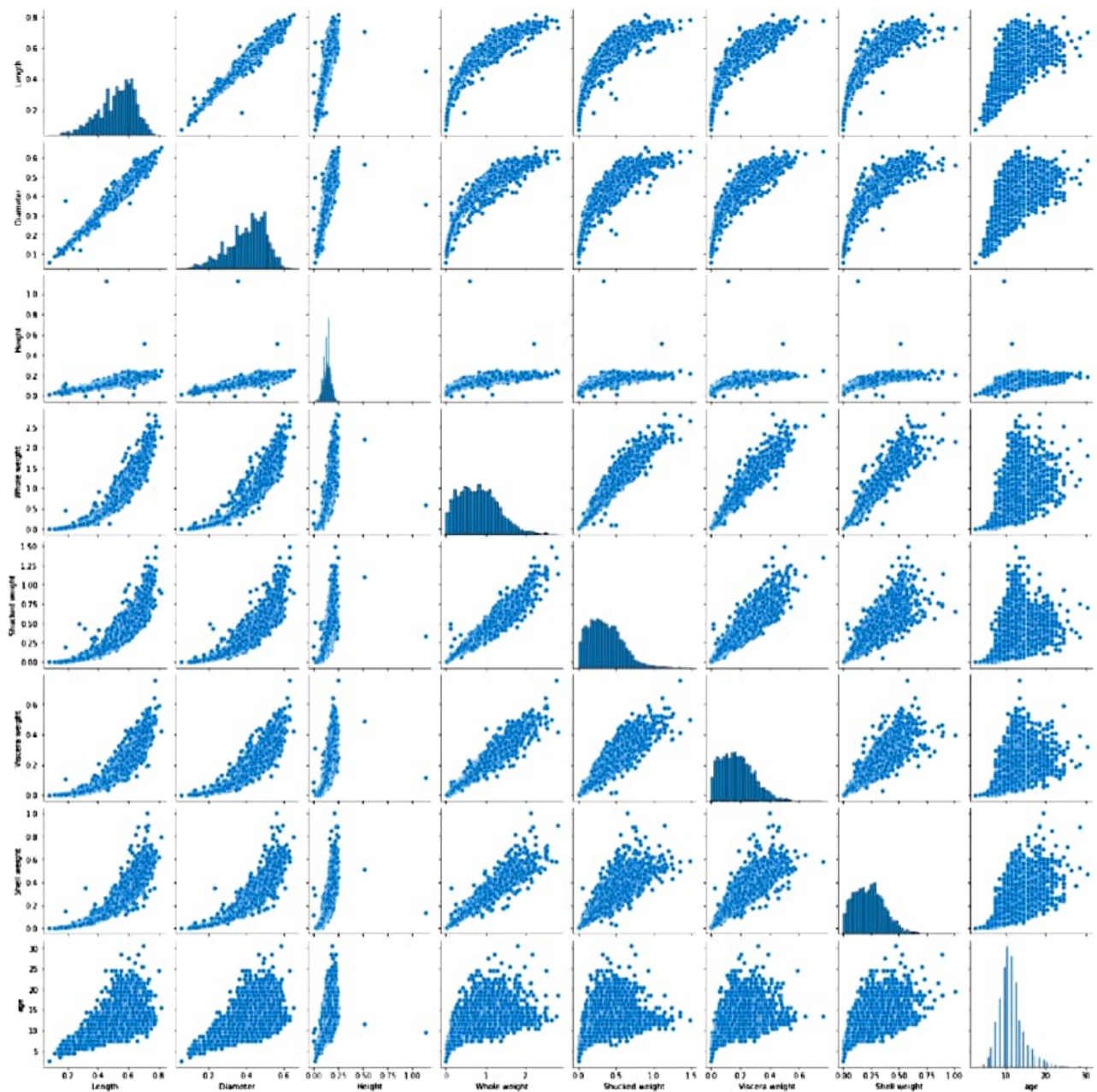
```
Out[5]:
```

|     | Length   | Diameter | Height   | Whole weight | Shucked weight | Viscera weight | Shell weight | age       |
|-----|----------|----------|----------|--------------|----------------|----------------|--------------|-----------|
| Sex |          |          |          |              |                |                |              |           |
| I   | 0.427746 | 0.326494 | 0.107996 | 0.431363     | 0.191035       | 0.092010       | 0.128182     | 9.390462  |
| M   | 0.561391 | 0.439287 | 0.151381 | 0.991459     | 0.432946       | 0.215545       | 0.281969     | 12.205497 |
| F   | 0.579093 | 0.454732 | 0.158011 | 1.046532     | 0.446188       | 0.230689       | 0.302010     | 12.629304 |

### 3. BIVARIATE ANALYSIS & MULTIVARIATE ANALYSIS

```
In [6]: numerical_features = df.select_dtypes(include = [np.number]).columns
sns.pairplot(df[numerical_features])
```

Out[6]:



### 4. Descriptive statistics

```
In [7]: df.describe()
```

Out[7]:

|       | Length      | Diameter    | Height      | Whole weight | Shucked weight | Viscera weight | Shell weight | age         |
|-------|-------------|-------------|-------------|--------------|----------------|----------------|--------------|-------------|
| 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     | 11.433684   |
| 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     | 2.500000    |
| 25%   | 0.450000    | 0.350000    | 0.115000    | 0.441500     | 0.186000       | 0.093500       | 0.130000     | 9.500000    |
| 50%   | 0.545000    | 0.425000    | 0.140000    | 0.799500     | 0.336000       | 0.171000       | 0.234000     | 10.500000   |
| 75%   | 0.615000    | 0.480000    | 0.165000    | 1.153000     | 0.502000       | 0.253000       | 0.329000     | 12.500000   |
| max   | 0.815000    | 0.650000    | 1.130000    | 2.825500     | 1.488000       | 0.760000       | 1.005000     | 30.500000   |

## 5. Check for Missing Values

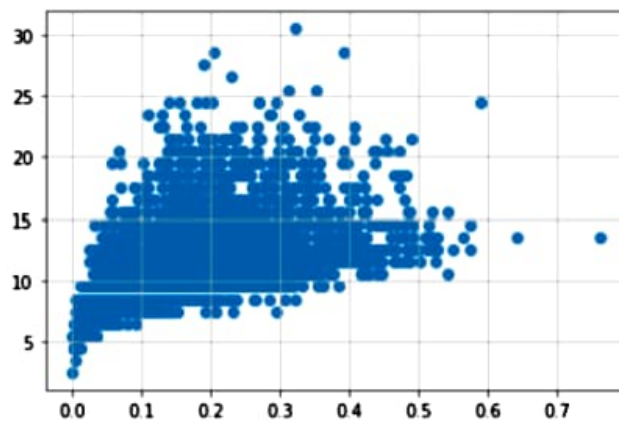
```
In [8]: df.isnull().sum()
```

```
Out[8]: Sex                0
Length                0
Diameter              0
Height               0
Whole weight          0
Shucked weight        0
Viscera weight         0
Shell weight          0
age                  0
dtype: int64
```

## 6. OUTLIER HANDLING

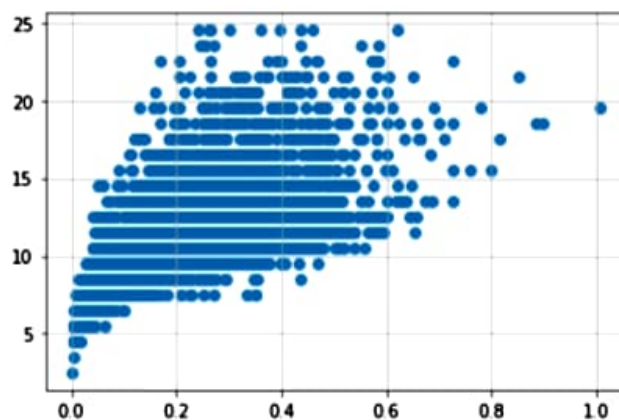
```
In [9]: df = pd.get_dummies(df)
dummy_data = df.copy()
```

```
In [10]: var = 'Viscera weight'
plt.scatter(x = df[var], y = df['age'],)
plt.grid(True)
```



```
In [11]: # outliers removal
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)
```

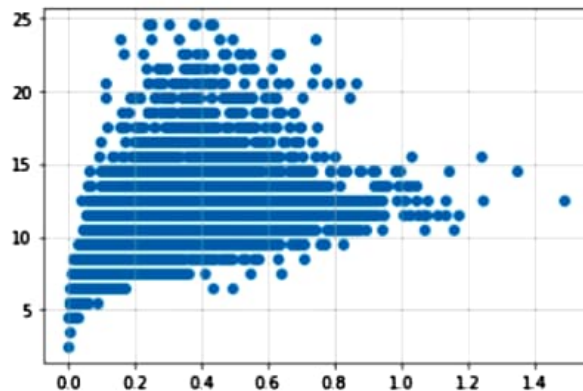
```
In [12]: var = 'Shell weight'
plt.scatter(x = df[var], y = df['age'],)
plt.grid(True)
#Outliers removal
df.drop(df[(df['Shell weight'] > 0.6) & (df['age'] < 25)].index, inplace=True)
df.drop(df[(df['Shell weight'] < 0.8) & (df['age'] > 25)].index, inplace=True)
```



In [13]:

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

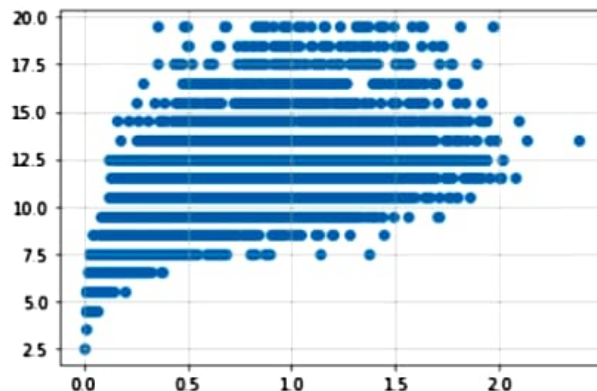
#Outlier removal
df.drop(df[(df['Shucked weight'] >= 1) & (df['age'] < 20)].index, inplace=True)
df.drop(df[(df['Shucked weight'] < 1) & (df['age'] > 20)].index, inplace=True)
```



In [14]:

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

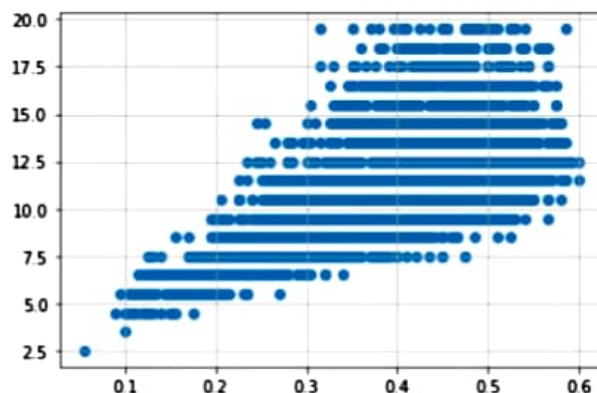
df.drop(df[(df['Whole weight'] >= 2.5) &
           (df['age'] < 25)].index, inplace = True)
df.drop(df[(df['Whole weight'] < 2.5) &
           (df['age'] > 25)].index, inplace = True)
```



In [15]:

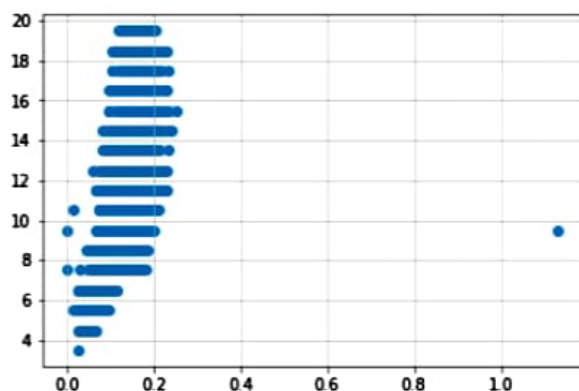
```
var = 'Diameter'
plt.scatter(x = df[var], y = df['age'])
plt.grid(True)

df.drop(df[(df['Diameter'] < 0.1) &
           (df['age'] < 5)].index, inplace = True)
df.drop(df[(df['Diameter'] < 0.6) &
           (df['age'] > 25)].index, inplace = True)
df.drop(df[(df['Diameter'] >= 0.6) &
           (df['age'] < 25)].index, inplace = True)
```



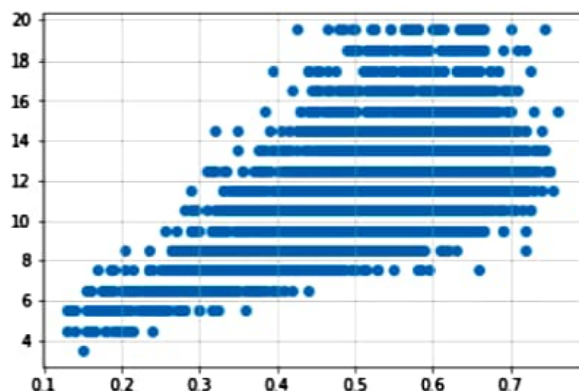


```
In [16]: var = 'Height'
plt.scatter(x = df[var], y = df['age'])
plt.grid(True)
df.drop(df[(df['Height'] > 0.4) &
          (df['age'] < 15)].index, inplace = True)
df.drop(df[(df['Height'] < 0.4) & (
df['age'] > 25)].index, inplace = True)
```



```
In [17]: var = 'Length'
plt.scatter(x = df[var], y = df['age'])
plt.grid(True)

df.drop(df[(df['Length'] < 0.1) &
          (df['age'] < 5)].index, inplace = True)
df.drop(df[(df['Length'] < 0.8) & (
df['age'] > 25)].index, inplace = True)
df.drop(df[(df['Length'] >= 0.8) & (
df['age'] < 25)].index, inplace = True)
```



## 7. Categorical columns

```
In [18]: numerical_features = df.select_dtypes(include = [np.number]).columns
categorical_features = df.select_dtypes(include = [np.object]).columns
```

/usr/local/lib/python3.7/dist-packages/ipykernel\_launcher.py:2: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warning, use `object` by itself. Doing this will not modify any behavior and is safe.  
Deprecated in NumPy 1.20; for more details and guidance: <https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations>

```
In [19]: numerical_features
```

```
Out[19]: Index(['Length', 'Diameter', 'Height', 'Whole weight', 'Shucked weight',
              'Viscera weight', 'Shell weight', 'age', 'Sex_F', 'Sex_I', 'Sex_M'],
              dtype='object')
```

```
In [20]: categorical_features
```

```
Out[20]: Index([], dtype='object')
```

## ENCODING

```
In [21]: from sklearn.preprocessing import LabelEncoder  
le=LabelEncoder()  
print(df.Length.value_counts())
```

```
0.575    93  
0.625    91  
0.580    89  
0.550    89  
0.620    83  
..  
0.220     2  
0.150     1  
0.755     1  
0.135     1  
0.760     1  
Name: Length, Length: 126, dtype: int64
```

## 8. Split the dependent and independent variables

```
In [22]: x=df.iloc[:,5]  
x
```

```
Out[22]:
```

|      | Length | Diameter | Height | Whole weight | Shucked weight |
|------|--------|----------|--------|--------------|----------------|
| 0    | 0.455  | 0.365    | 0.095  | 0.5140       | 0.2245         |
| 1    | 0.350  | 0.265    | 0.090  | 0.2255       | 0.0995         |
| 2    | 0.530  | 0.420    | 0.135  | 0.6770       | 0.2565         |
| 3    | 0.440  | 0.365    | 0.125  | 0.5160       | 0.2155         |
| 4    | 0.330  | 0.255    | 0.080  | 0.2050       | 0.0895         |
| ...  | ...    | ...      | ...    | ...          | ...            |
| 4172 | 0.565  | 0.450    | 0.165  | 0.8870       | 0.3700         |
| 4173 | 0.590  | 0.440    | 0.135  | 0.9660       | 0.4390         |
| 4174 | 0.600  | 0.475    | 0.205  | 1.1760       | 0.5255         |
| 4175 | 0.625  | 0.485    | 0.150  | 1.0945       | 0.5310         |
| 4176 | 0.710  | 0.555    | 0.195  | 1.9485       | 0.9455         |

3995 rows × 5 columns

```
In [23]: y=df.iloc[:,5:]  
y
```

```
Out[23]:
```

|      | Viscera weight | Shell weight | age  | Sex_F | Sex_I | Sex_M |
|------|----------------|--------------|------|-------|-------|-------|
| 0    | 0.1010         | 0.1500       | 16.5 | 0     | 0     | 1     |
| 1    | 0.0485         | 0.0700       | 8.5  | 0     | 0     | 1     |
| 2    | 0.1415         | 0.2100       | 10.5 | 1     | 0     | 0     |
| 3    | 0.1140         | 0.1550       | 11.5 | 0     | 0     | 1     |
| 4    | 0.0395         | 0.0550       | 8.5  | 0     | 1     | 0     |
| ...  | ...            | ...          | ...  | ...   | ...   | ...   |
| 4172 | 0.2390         | 0.2490       | 12.5 | 1     | 0     | 0     |
| 4173 | 0.2145         | 0.2605       | 11.5 | 0     | 0     | 1     |
| 4174 | 0.2875         | 0.3080       | 10.5 | 0     | 0     | 1     |
| 4175 | 0.2610         | 0.2960       | 11.5 | 1     | 0     | 0     |
| 4176 | 0.3765         | 0.4950       | 13.5 | 0     | 0     | 1     |

3995 rows × 6 columns

## 9. Feature Scaling

```
In [24]: from sklearn.preprocessing import StandardScaler
         ss=StandardScaler()
         x_train=ss.fit_transform(x_train)
```

```
In [25]: mlrpred=mlr.predict(x_test[0:9])
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/base.py:444: UserWarning: X has feature names, but LinearRegression was fitted without feature names
  f"X has feature names, but {self.__class__.__name__} was fitted without"
```

```
In [26]: mlrpred
```

```
Out[26]: array([[ 0.23339315,  0.30675115, 12.64851662,  0.41303667,  0.21495648,
                  0.37200685],
                [ 0.29781617,  0.38717341, 13.7465214 ,  0.49950512,  0.07794454,
                  0.42255034],
                [ 0.31212505,  0.40571258, 14.08610548,  0.51832591,  0.0551311 ,
                  0.42654299],
                [ 0.28929529,  0.37330035, 13.37905462,  0.47977802,  0.10981522,
                  0.41040676],
                [ 0.27398024,  0.35592095, 13.33560842,  0.45982985,  0.14616275,
                  0.39400739],
                [ 0.26092694,  0.34070882, 13.08733159,  0.44915929,  0.1554206 ,
                  0.3954201 ],
                [ 0.28798179,  0.37424621, 13.57347267,  0.48251798,  0.10441507,
                  0.41306695],
                [ 0.35187797,  0.45729677, 14.8824873 ,  0.57548893, -0.03481323,
                  0.45932431],
                [ 0.19358839,  0.2544071 , 11.75629202,  0.35124826,  0.30061553,
                  0.34813621]])
```

## 10. Train , Test , Split

```
In [27]: from sklearn.model_selection import train_test_split
         x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2)
```

## 11. Model building

```
In [28]: from sklearn.linear_model import LinearRegression
         mlr=LinearRegression()
         mlr.fit(x_train,y_train)
```

```
Out[28]: LinearRegression()
```

## 12 & 13. Train and Test the model

```
In [29]: x_test[0:5]
```

```
Out[29]:
```

|      | Length | Diameter | Height | Whole weight | Shucked weight |
|------|--------|----------|--------|--------------|----------------|
| 144  | 0.475  | 0.375    | 0.130  | 0.5175       | 0.2075         |
| 2783 | 0.600  | 0.495    | 0.185  | 1.1145       | 0.5055         |
| 4139 | 0.635  | 0.495    | 0.175  | 1.2355       | 0.5205         |
| 2801 | 0.640  | 0.515    | 0.080  | 1.0420       | 0.5150         |
| 3185 | 0.590  | 0.415    | 0.150  | 0.8805       | 0.3645         |

```
In [30]: y_test[0:5]
```

```
Out[30]:
```

|      | Viscera weight | Shell weight | age  | Sex_F | Sex_I | Sex_M |
|------|----------------|--------------|------|-------|-------|-------|
| 144  | 0.1165         | 0.170        | 11.5 | 0     | 0     | 1     |
| 2783 | 0.2635         | 0.367        | 12.5 | 0     | 0     | 1     |
| 4139 | 0.3085         | 0.347        | 11.5 | 1     | 0     | 0     |
| 2801 | 0.1755         | 0.175        | 11.5 | 0     | 0     | 1     |
| 3185 | 0.2340         | 0.235        | 12.5 | 1     | 0     | 0     |

## 14. Measure the performance using metrics

```
In [31]: from sklearn.metrics import r2_score  
r2_score(mlr.predict(x_test), y_test)
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/base.py:444: UserWarning: X has feature names, but LinearRegression was fitted without feature names  
  f"X has feature names, but {self.__class__.__name__} was fitted without"
```

```
Out[31]: -48.09199504251874
```