

## Assignment-4

<b>Date</b>	<b>10.11.2022</b>
<b>Team Id</b>	<b>PNT2022TMID24138</b>
<b>Project Name</b>	<b>Estimated The Crop Yield using Data Analytics</b>
<b>Team Leader</b>	<b>D. Sandhiya</b>
<b>Team Member</b>	<b>M. Ruthra P.Priyadharshini D.Tejasvani</b>

```

IMPORTING LIBRARIES

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

from google.colab import drive
drive.mount('/content/drive')

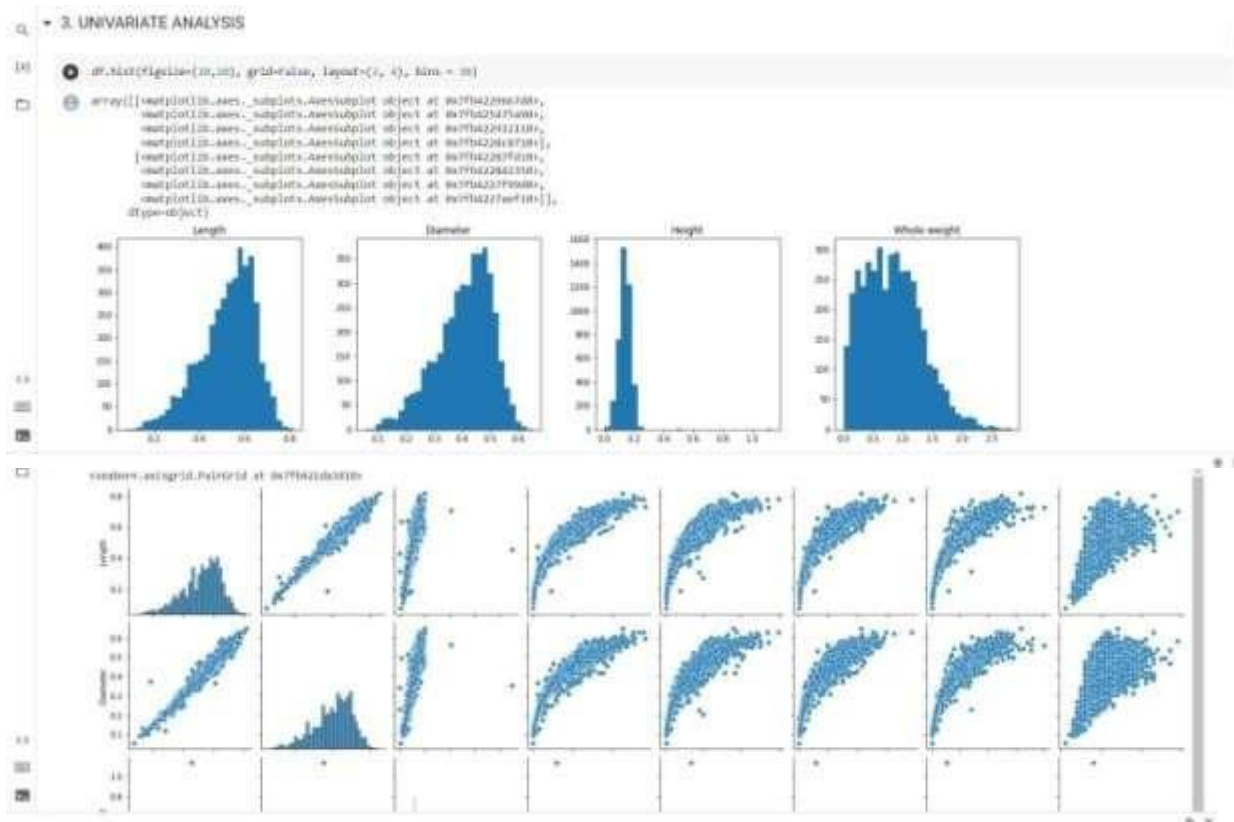
!df=pd.read_csv('/content/drive/MyDrive/sem assignment 4/walrus.csv')

df.size

17503

3. UNIVARIATE ANALYSIS

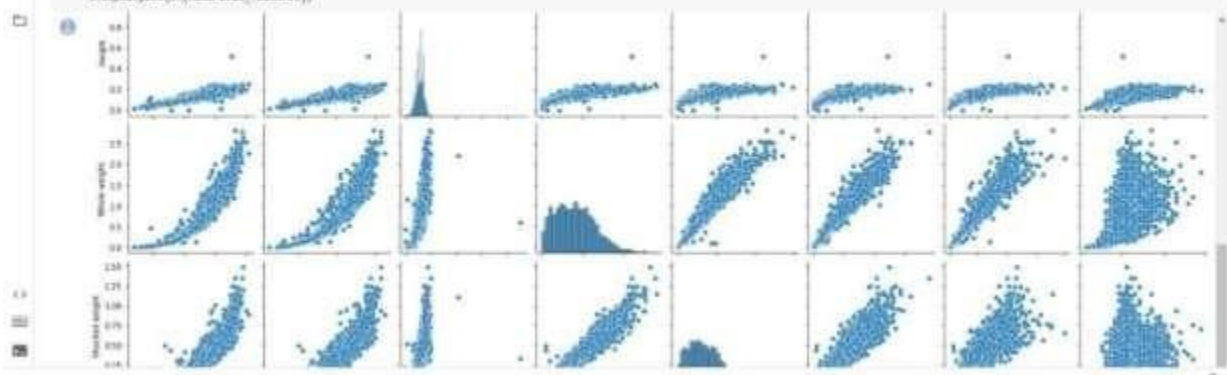
```





### 3. BIVARIATE ANALYSIS & MULTIVARIATE ANALYSIS

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



### 4. Descriptive statistics

```
df.describe()
```

	length	diameter	weight	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	0.523962	0.427881	0.132018	0.528742	0.280367	0.182284	0.238831	9.533664
std	0.120080	0.095240	0.041627	0.480388	0.221903	0.130814	0.136025	3.224169
min	0.075000	0.050000	0.000000	0.002000	0.001000	0.000000	0.001000	1.000000
25%	0.400000	0.300000	0.115000	0.441500	0.180000	0.092000	0.130000	8.000000
50%	0.545000	0.425000	0.140000	0.796000	0.380000	0.171000	0.234000	9.500000
75%	0.618000	0.480000	0.180000	1.150000	0.602000	0.253000	0.329000	11.000000
max	0.810000	0.650000	1.130000	2.825000	1.480000	0.780000	1.050000	20.000000

### 5. Check for Missing Values

```
df.isnull().sum()
```

```
length    0
diameter  0
weight    0
whole weight    0
shucked weight    0
viscera weight  0
shell weight  0
rings      0
dtype: int64
```

### 6. OUTLIER HANDLING

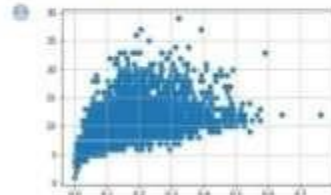
```
df = pd.get_dummies(df)
dummy_data = df.copy()

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

## 6. OUTLIER HANDLING

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

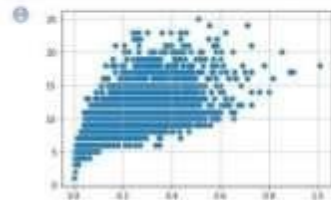
```
var = 'viscous weight'
plt.scatter(x=df[var], y=df['kings'],)
plt.grid(True)
```



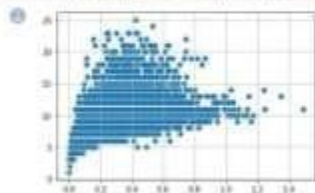
```
[ ] # outliers removal
df.drop(df[(df['viscous weight'] > 0.7) & (df['kings'] < 20)].index, inplace=True)
```

```
[ ] # outliers removal
df.drop(df[(df['viscous weight'] > 0.7) & (df['kings'] < 20)].index, inplace=True)
df.drop(df[(df['viscous weight'] < 0.1) & (df['kings'] > 25)].index, inplace=True)
```

```
var = 'shell weight'
plt.scatter(x=df[var], y=df['kings'],)
plt.grid(True)
# outliers removal
df.drop(df[(df['shell weight'] > 0.6) & (df['kings'] < 20)].index, inplace=True)
df.drop(df[(df['shell weight'] < 0.1) & (df['kings'] > 25)].index, inplace=True)
```



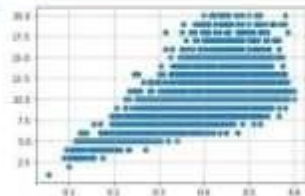
```
var = 'shock weight'
plt.scatter(x=df[var], y=df['kings'],)
plt.grid(True)
# outliers removal
df.drop(df[(df['shock weight'] > 1) & (df['kings'] < 20)].index, inplace=True)
df.drop(df[(df['shock weight'] < 0.1) & (df['kings'] > 25)].index, inplace=True)
```



```
[ ] var = 'hole weight'
plt.scatter(x=df[var], y=df['kings'],)
plt.grid(True)
df.drop(df[(df['hole weight'] > 2.5) & (df['kings'] < 20)].index, inplace=True)
```

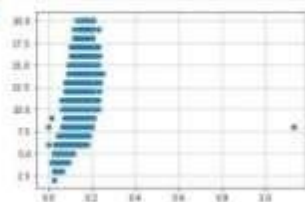
```
var = 'diameter'
plt.scatter(x = df[var], y = df['kings'])
plt.grid(True)

df.drop(df[(df['diameter'] < 0.1) &
          (df['kings'] < 5)].index, inplace = True)
df.drop(df[(df['diameter'] > 0.4) &
          (df['kings'] > 25)].index, inplace = True)
df.drop(df[(df['diameter'] > 0.4) &
          (df['kings'] < 5)].index, inplace = True)
```



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

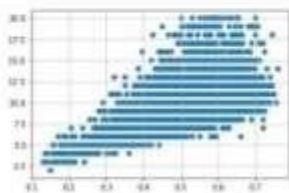
```
var = 'weight'
plt.scatter(x = df[var], y = df['kings'])
plt.grid(True)
df.drop(df[(df['weight'] > 0.4) &
          (df['kings'] < 5)].index, inplace = True)
df.drop(df[(df['weight'] > 0.4) &
          (df['kings'] > 25)].index, inplace = True)
```



```
} } var = 'length'
plt.scatter(x = df[var], y = df['kings'])
plt.grid(True)
```

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

df.drop(df[(df['length'] < 0.1) &
          (df['kings'] < 5)].index, inplace = True)
df.drop(df[(df['length'] > 0.4) &
          (df['kings'] > 25)].index, inplace = True)
df.drop(df[(df['length'] > 0.4) &
          (df['kings'] < 5)].index, inplace = True)
```



## 7. Categorical columns

```
[ ] 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/pymongo/_helpers.py:2: DeprecationWarning: 'np.object' is a deprecated alias for the builtin 'object'. To silence this warning, use 'object' by
deprecated in NumPy 1.18; for more details and guidance: https://numpy.org/doc/stable/1.18-0-notes/deprecations.html
1

numerical_features
Index(['length', 'diameter', 'weight', 'shale weight', 'shucked weight',
       'viscera weight', 'shell weight', 'kings', 'sex_F', 'sex_M'],
      dtype=object)

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

## ENCODING

```
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
print(le.fit_transform(df.length.value_counts()))
```

```
0.575  93
0.180  92
0.550  90
0.625  84
0.680  85
...
0.110  3
0.120  0
0.715  5
0.100  0
0.700  5
Name: length, length: int64, dtype: int64
```

## 8. Split the dependent and independent variables

```
[ ] x=df.iloc[:,0:5]
y =

length diameter weight shale weight shucked weight
```

## 8. Split the dependent and independent variables

```
x=df.iloc[:,0:5]
y =
```

	length	diameter	weight	shale weight	shucked weight
0	0.455	0.365	0.085	0.0140	0.2245
1	0.380	0.265	0.080	0.0056	0.0000
2	0.530	0.420	0.105	0.0775	0.2565
3	0.440	0.365	0.125	0.0755	0.2150
4	0.330	0.205	0.080	0.0000	0.0885
...	...	...	...	...	...
4172	0.565	0.430	0.145	0.0875	0.3750
4173	0.590	0.440	0.135	0.0950	0.4300
4174	0.600	0.475	0.215	1.1700	0.5055
4175	0.625	0.485	0.190	1.0840	0.5210
4176	0.710	0.505	0.180	1.0485	0.5405

4049 rows × 5 columns

```

y=07.11ac[1:,1]
y

```

(1)

	Viscera weight	Shell weight	Rings	Sex_F	Sex_I	Sex_M
0	0.1010	0.1500	10	0	0	1
1	0.0480	0.0700	7	0	0	1
2	0.0415	0.2100	8	1	0	0
3	0.1140	0.1650	10	0	0	1
4	0.0380	0.0000	7	0	1	0
...	...	...	...	...	...	...
4172	0.2390	0.2480	11	1	0	0
4173	0.2140	0.2800	10	0	0	1
4174	0.2670	0.3080	9	0	0	1
4175	0.2010	0.0900	10	1	0	0
4176	0.3700	0.4850	12	0	0	1

4183 rows x 6 columns

```

9. Feature Scaling
from sklearn.preprocessing import StandardScaler
ss=StandardScaler()
ss.fit_train(x)

10. Train , Test, Split

from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)

11. Model building

from sklearn.linear_model import LinearRegression
slr=LinearRegression()
slr.fit(x_train,y_train)

LinearRegression()

```

```

12 & 13. Train and Test the model

x_test[0:]
array([[0.17386052, 0.10558061, 0.09524821, -0.18950158, 0.01702011],
       [0.16651809, 0.08157811, 0.07773846, 0.46723804, 0.25809106],
       [0.10227239, -0.0728862, -0.38923226, -0.19680227, -0.10157271],
       [0.10726123, -0.24728802, -0.49871389, -0.10884188, -0.11771341],
       [0.05648993, -0.55016289, -0.6745113, -1.48138971, -1.88133127]])

y_test[0:]

```

	Viscera weight	Shell weight	Rings	Sex_F	Sex_I	Sex_M
3484	0.1500	0.1840	8	0	1	0
2458	0.2520	0.3680	10	1	0	0
3252	0.1420	0.1700	12	0	0	1
2648	0.1200	0.1300	7	0	1	0
3638	0.0810	0.1050	8	0	1	0

#### 14. Measure the performance using metrics

```
from sklearn.metrics import r2_score  
r2_score(y_test, predict(x_test, y_test))
```

```
-1.0000000000000001
```