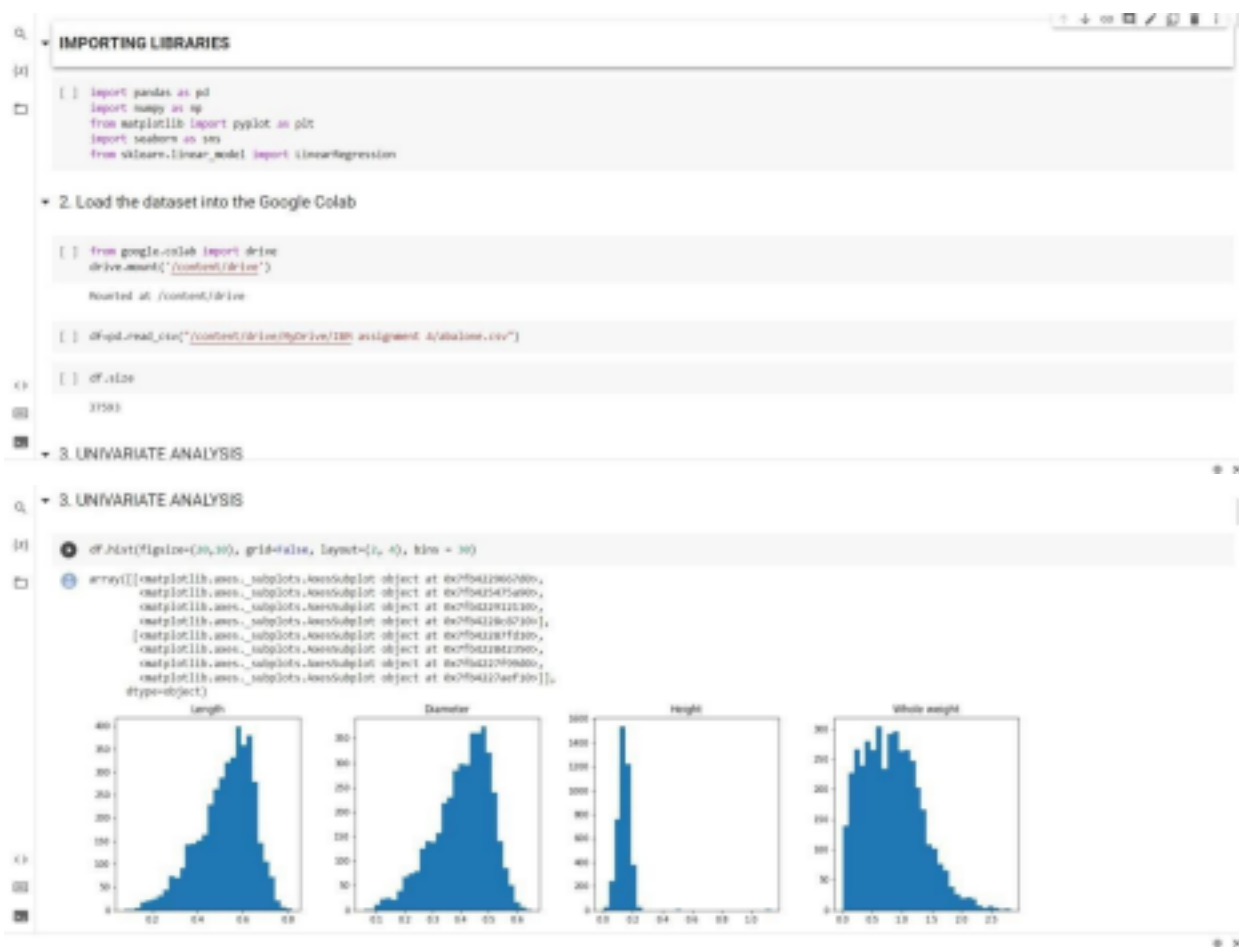


Assignment-4

Date	03.11.2022
Team Id	PNT2022TMID24195
Project Name	Estimated The Crop Yelid using Data Analytics
Team Leader	Divi Nandhu
Team Member	Bharathkumar S Induja H Chaithanya Babu K



3. BIVARIATE ANALYSIS & MULTIVARIATE ANALYSIS

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

seaborn.relgrid.FacetGrid at 0x7f424d4d0000

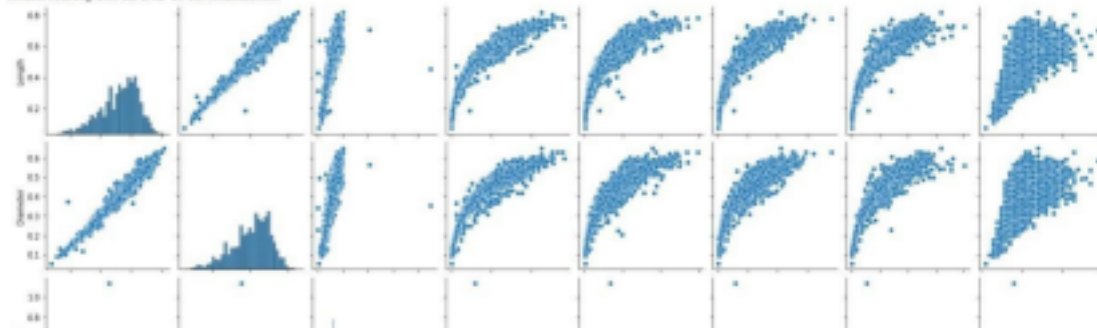
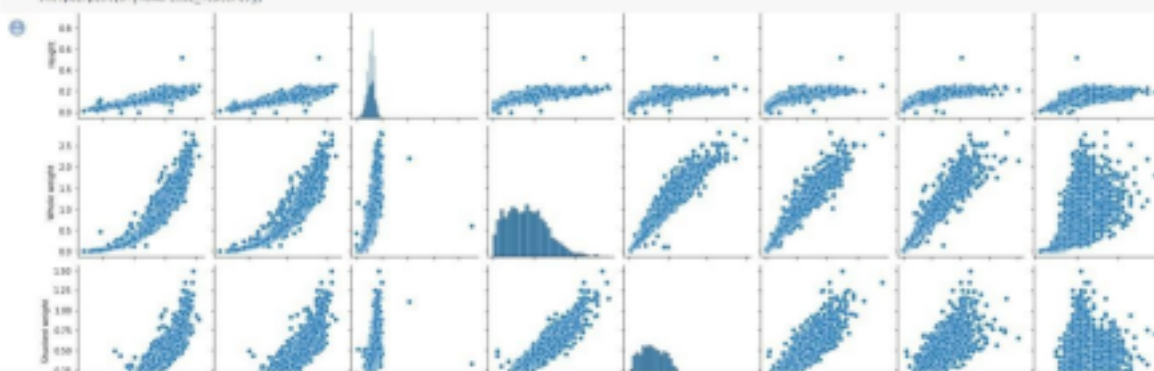


Figure 1

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	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	0.529902	0.407901	0.136518	0.828742	0.358567	0.108594	0.238001	8.003684
std	0.129003	0.098048	0.047827	0.486389	0.227983	0.108614	0.136293	3.224189
min	0.075000	0.058000	0.008000	0.002000	0.001000	0.008000	0.001500	1.000000
25%	0.450000	0.300000	0.115000	0.441500	0.156000	0.009500	0.130000	8.000000
50%	0.545000	0.425000	0.140000	0.798500	0.338000	0.171000	0.234000	8.000000
75%	0.615000	0.450000	0.155000	1.153000	0.502000	0.253000	0.329000	11.000000
max	0.815000	0.650000	1.130000	2.825000	1.488000	0.768000	1.008000	28.000000

5. Check for Missing Values

```
[ ] df.isna().sum()
```

```
Size      0
Length    0
Diameter  0
Height    0
Whole weight  0
Shucked weight  0
Viscera weight  0
Shell weight  0
Rings     0
dtypes: int64
```

6. OUTLIER HANDLING

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

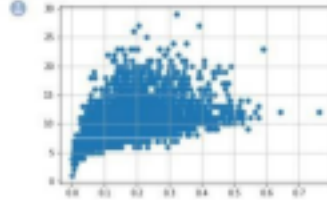
[ ] var = 'Viscera weight'
plt.scatter(x=df[var], y=df['Rings'],)
plt.grid(True)
```

Figure 2

6. OUTLIER HANDLING

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

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

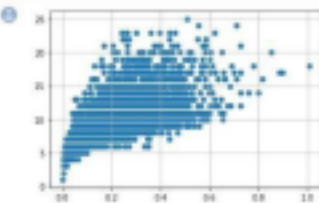


```
[ ] # outliers removal
df.drop(df[(df['viscera weight'] > 0.5) & (df['kings'] < 20)].index, inplace=True)
```

```
[ ] # outliers removal
df.drop(df[(df['viscera weight'] > 0.5) & (df['kings'] < 20)].index, inplace=True)
df.drop(df[(df['viscera weight'] < 0.5) & (df['kings'] > 25)].index, inplace=True)
```

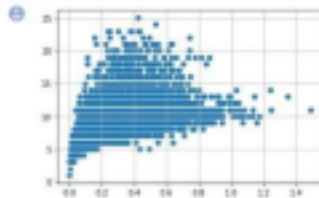
```
var = 'shell weight'
plt.scatter(x = df[var], y = df['kings'],)
plt.grid(True)

# outliers removal
df.drop(df[(df['shell weight'] > 0.5) & (df['kings'] < 20)].index, inplace=True)
df.drop(df[(df['shell weight'] < 0.5) & (df['kings'] > 25)].index, inplace=True)
```



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

# outliers removal
df.drop(df[(df['shucked weight'] > 1) & (df['kings'] < 20)].index, inplace=True)
df.drop(df[(df['shucked weight'] < 1) & (df['kings'] > 25)].index, inplace=True)
```

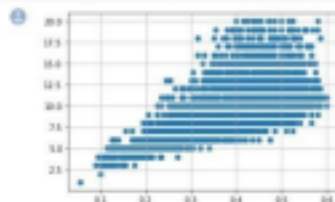


```
[ ] var = 'whole weight'
plt.scatter(x = df[var], y = df['kings'],)
plt.grid(True)

df.drop(df[(df['whole weight'] > 1.5) &
(df['kings'] < 20)].index, inplace = True)
```

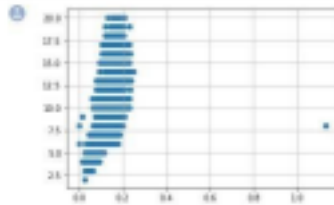
```
var = 'diameter'
plt.scatter(x = df[var], y = df['kings'],)
plt.grid(True)

df.drop(df[(df['diameter'] > 0.1) &
(df['kings'] < 5)].index, inplace = True)
df.drop(df[(df['diameter'] < 0.1) &
(df['kings'] > 25)].index, inplace = True)
df.drop(df[(df['diameter'] > 0.1) &
(df['kings'] < 20)].index, inplace = True)
```



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

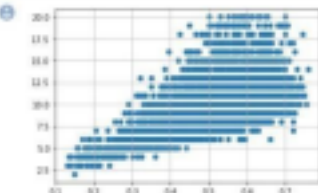
```
var = 'weight'
plt.scatter(x = df[var], y = df['kings'])
plt.grid(True)
df.drop(df[(df['weight'] > 0.4) &
          (df['kings'] < 1)], index, inplace = True)
df.drop(df[(df['weight'] < 0.4) &
          (df['kings'] > 1)], index, inplace = True)
```



```
[ ] 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'] < 1)], index, inplace = True)
df.drop(df[(df['length'] > 0.8) &
          (df['kings'] > 1)], index, inplace = True)
df.drop(df[(df['length'] > 0.8) &
          (df['kings'] < 1)], 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/ipykernel_launcher.py:1: DeprecationWarning: 'np.object' is a deprecated alias for the builtin 'object'. To silence this warning, use 'object' by deprecated in NumPy 1.20; for more details and guidance: <https://numpy.org/doc/stable/release/1.20.0-notes/3d0c46d2deprecation.html>

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

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

ENCODING

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

```
0.575    95
0.580    92
0.590    92
0.625    91
0.680    85
...
0.125     1
0.130     1
0.755     1
0.150     1
0.760     1
Name: length, length: 126, dtype: object
```

8. Split the dependent and independent variables

```
X = df.iloc[:,15]
y =
```

length diameter weight shale weight shucked weight

8. Split the dependent and independent variables

add_slice(i, n)

x

	length	diameter	weight	whole weight	shucked weight
0	8.438	0.365	0.085	0.8143	0.2348
1	8.586	0.265	0.080	0.2255	0.8895
2	8.530	0.420	0.125	0.8773	0.2565
3	8.448	0.365	0.125	0.9180	0.2158
4	8.530	0.255	0.080	0.2050	0.8895
...
4172	8.588	0.480	0.185	0.8872	0.3708
4173	8.586	0.440	0.135	0.8690	0.4390
4174	8.636	0.475	0.295	1.1790	0.5255
4175	8.628	0.485	0.180	1.2948	0.8310
4176	8.710	0.555	0.185	1.8455	0.9455

4049 rows × 6 columns



