

Library Import

```
In [ ]: import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns
```

SK-Learn Library Import

```
In [ ]: ''' SK-Learn Library Import'''
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import RandomizedLasso, LassoLarsCV
from sklearn.exceptions import ConvergenceWarning
from sklearn.model_selection import train_test_split
from sklearn.metrics import make_scorer, accuracy_score, confusion_matrix
import sklearn.datasets
```

```
In [ ]: %scipy, Stats Library'''
from scipy.stats import skew
```

```
In [ ]: import warnings
warnings.filterwarnings('ignore')
```

```
In [ ]: sns.set()
%matplotlib inline
```

Display The head -> To Check if Data is Properly Imported

```
In [11]: df = pd.read_csv(r'https://raw.githubusercontent.com/drsrscientist/dataset1/master/abalone.csv')
df.head()
```

```
Out[11]:
```

	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

Feature Information of the DataSet

```
In [12]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4177 entries, 0 to 4176
Data columns (total 9 columns):
#   Column          Non-Null Count  Dtype
---  --
0    Sex              4177 non-null   object
1    Length           4177 non-null   float64
2    Diameter         4177 non-null   float64
3    Height           4177 non-null   float64
4    Whole weight     4177 non-null   float64
5    Shucked weight   4177 non-null   float64
6    Viscera weight   4177 non-null   float64
7    Shell weight     4177 non-null   float64
8    Rings            4177 non-null   int64
dtypes: float64(7), int64(1), object(1)
memory usage: 293.8+ KB

According to the Information:
1)No-Null data 2)1 - Object Type 3)7 - Float Type 4)1 - Int Type
Feature Distribution of data for Float and Int Data Type
```

```
In [13]: df.describe()
```

```
Out[13]:
```

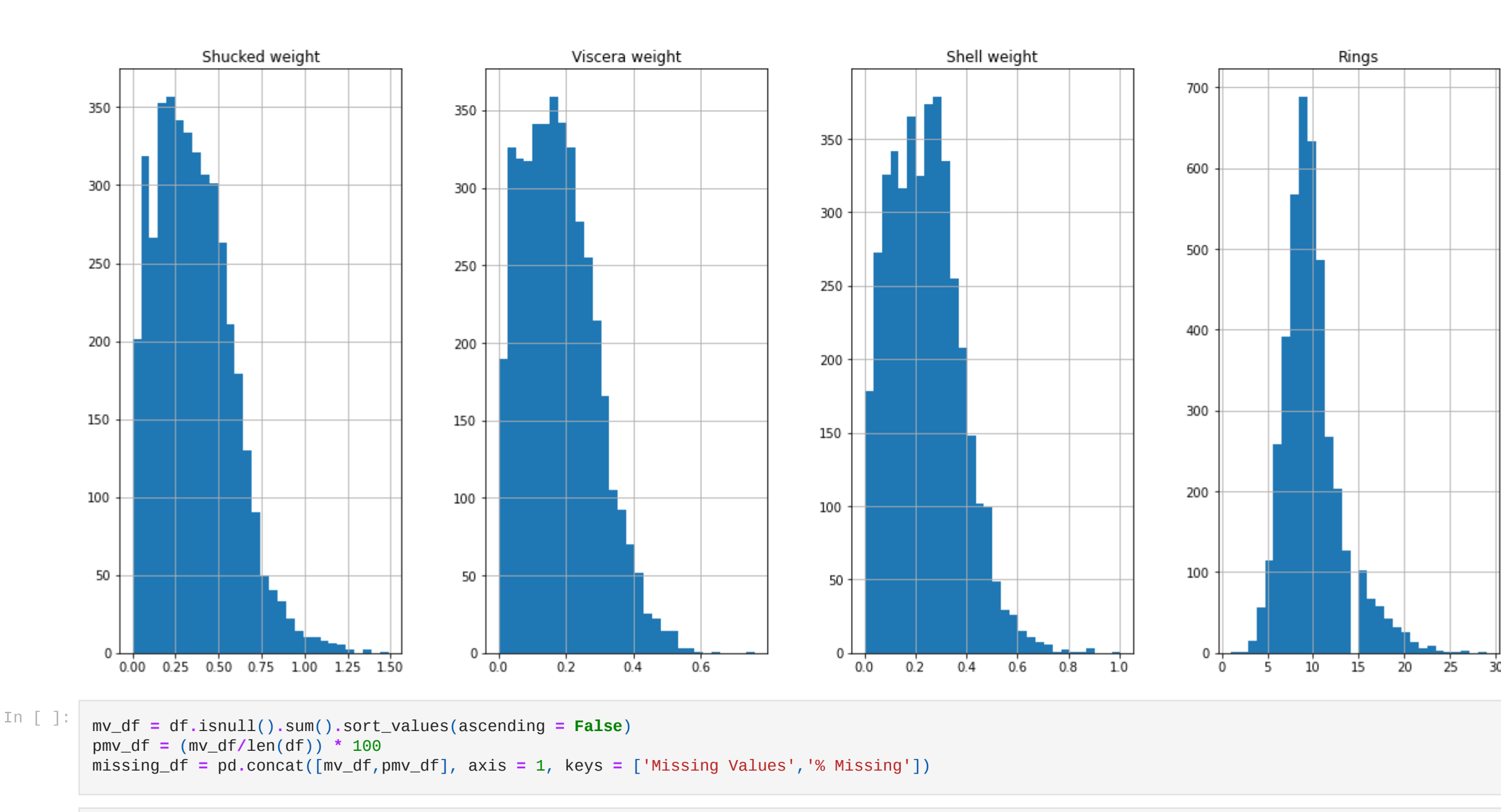
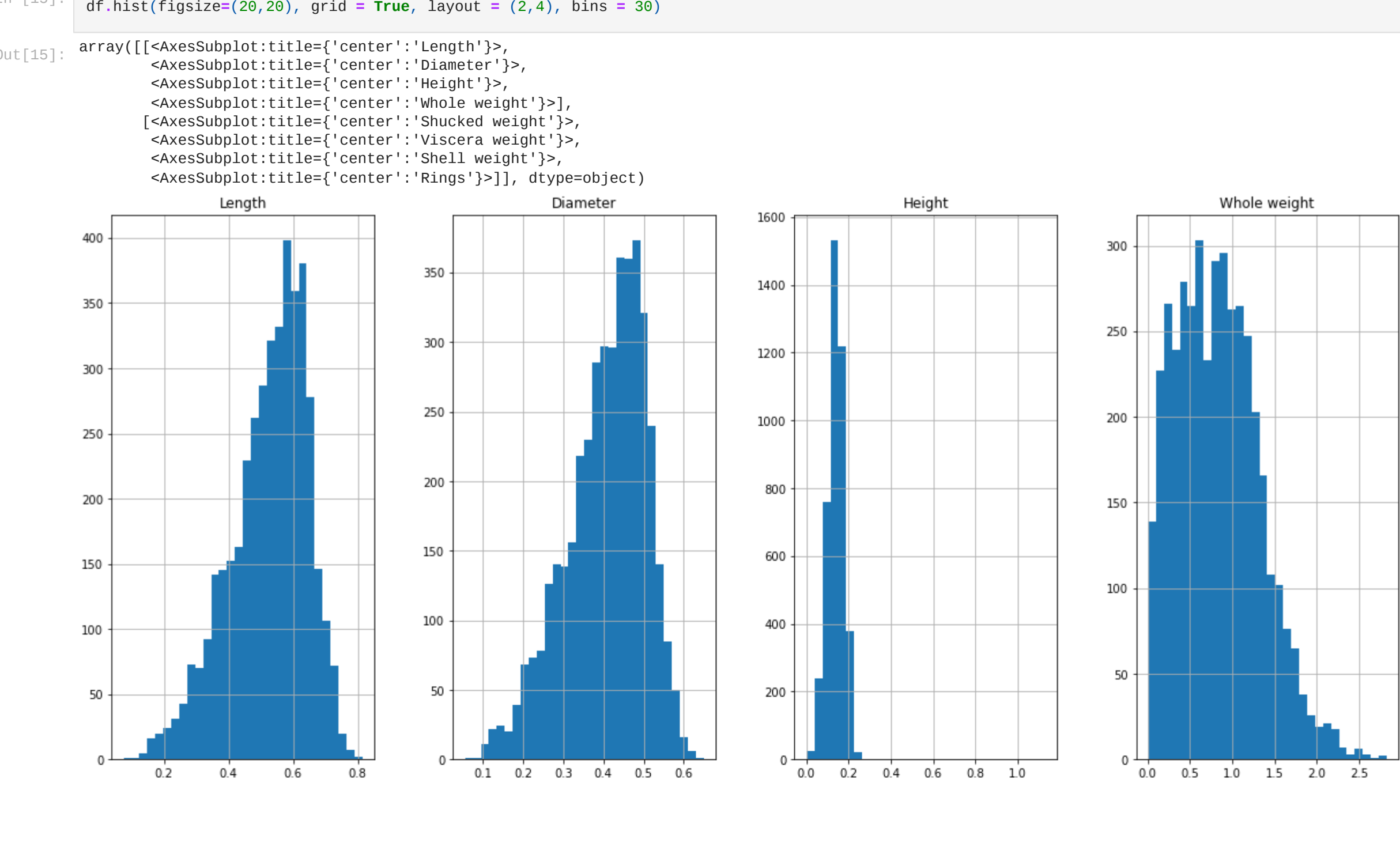
	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	Rings
count	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000
mean	0.523992	0.407881	0.139516	0.828742	0.359367	0.180594	0.238831	9.933684
std	0.120093	0.099240	0.041827	0.490389	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

According to Described Information:

```
In [ ]: 1)No Feature has Minimum Value = 0, except Height
2)All Features are Normally Distributed, ( Theoretically if feature is normally distributed, Mean = Median = Mode ).
3)But Features are close to Normality
4)All numerical, Except Sex
5)Each Feature has Different Scale
```

```
In [ ]: nf = df.select_dtypes(include=[np.number]).columns
cf = df.select_dtypes(include=[np.object]).columns
```

```
In [15]: df.hist(figsize=(20,20), grid = True, layout = (2,4), bins = 30)
```



```
In [ ]: mv_df = df.isnull().sum().sort_values(ascending = False)
pmv_df = (mv_df/len(df)) * 100
missing_df = pd.concat([mv_df,pmv_df], axis = 1, keys = ['Missing Values','% Missing'])
```

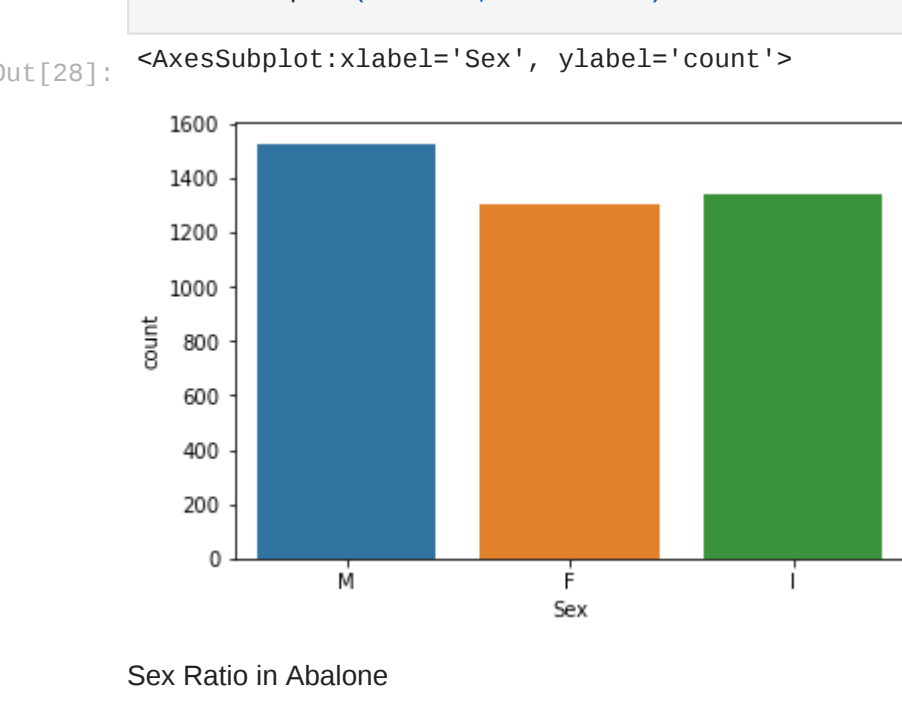
```
In [24]: print("Value Count of Rings Column")
print(df.Rings.value_counts())
print("\nPercentage of Rings Column")
print(df.Rings.value_counts(normalize = True))
```

```
Value Count of Rings Column
9      689
10     634
8      568
11     487
7      391
12     267
6      259
13     203
14     126
5      115
15     103
16      67
17      58
4       57
18     42
19     32
20     26
3       15
21     14
23      9
22      6
27      2
24      2
1        1
26      1
29      1
2        1
25      1
Name: Rings, dtype: int64

Percentage of Rings Column
9      0.164951
10     0.151784
8      0.135903
11     0.116591
7      0.093608
6      0.063921
5      0.062086
13     0.048599
14     0.038165
4      0.027532
15     0.024659
16     0.018048
17     0.013886
3      0.013646
18     0.010955
19     0.007661
20     0.006225
9      0.003591
21     0.003352
23     0.002155
22     0.001436
27     0.000479
24     0.000479
1      0.000239
26     0.000239
29     0.000239
2      0.000239
25     0.000239
Name: Rings, dtype: float64
```

Visualization

```
In [28]: import seaborn as sns
sns.countplot(x='Sex', data = df)
```



Sex Ratio in Abalone

```
In [29]: print("\nSex Count in Percentage")
print(df.Sex.value_counts(normalize = True))
print("\nSex Count in Numbers")
print(df.Sex.value_counts())
```

```
Sex Count in Percentage
M      0.365813
I      0.321283
F      0.312904
Name: Sex, dtype: float64

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

Small Feature Engineering, Deriving Age from Rings Column, Age = Rings + 1.5"

```
In [30]: df['Age'] = df['Rings'] + 1.5
df['Age'].head(5)
```

```
Out[30]:
```

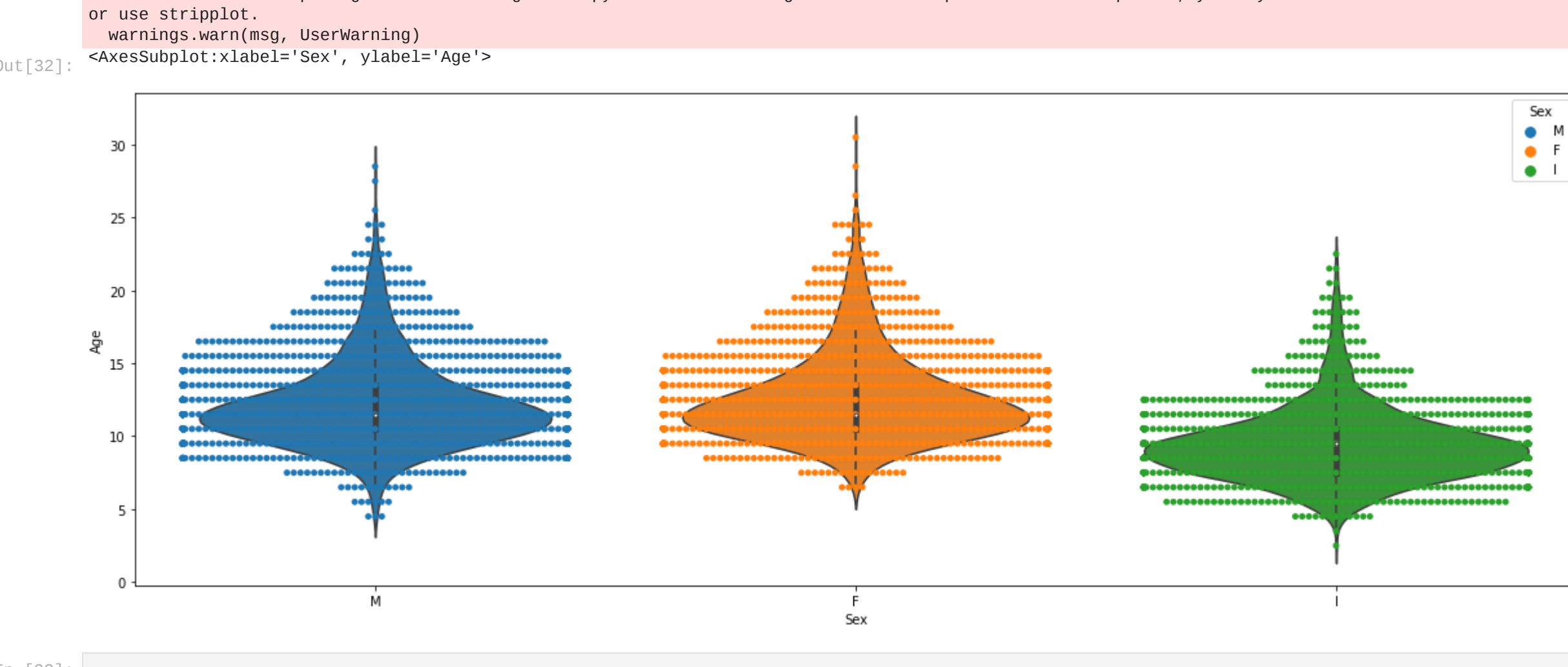
0	16.5
1	8.5
2	10.5
3	11.5
4	8.5

Name: Age, dtype: float64

"Sex and Age Visualization"

```
In [32]: import seaborn as sns
plt.figure(figsize = (20,7))
sns.swarmplot(x = 'Sex', y = 'Age', data = df, hue = 'Sex')
sns.wiplot(x = 'Sex', y = 'Age', data = df)
```

D:\ANACONDA\lib\site-packages\seaborn\categorical.py:1296: UserWarning: 56.2% of the points cannot be placed; you may want to decrease the size of the markers or use stripplot.
warnings.warn(msg, UserWarning)
D:\ANACONDA\lib\site-packages\seaborn\categorical.py:1296: UserWarning: 52.2% of the points cannot be placed; you may want to decrease the size of the markers or use stripplot.
warnings.warn(msg, UserWarning)
D:\ANACONDA\lib\site-packages\seaborn\categorical.py:1296: UserWarning: 58.5% of the points cannot be placed; you may want to decrease the size of the markers or use stripplot.
warnings.warn(msg, UserWarning)



```
In [33]: df.groupby('Sex')[['Length', 'Diameter', 'Height', 'Whole weight', 'Shucked weight', 'Viscera weight', 'Shell weight', 'Age']].mean().sort_values(by = 'Age',ascending = False)
```

```
Out[33]:
```

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

Preprocessing Data for the Model

```
In [ ]: df['Sex'] = LabelEncoder().fit_transform(df['Sex'].tolist())
```

```
In [ ]: transformed_sex_feature = OneHotEncoder().fit_transform(df['Sex'].values.reshape(-1,1)).toarray()
df_sex_encoded = pd.DataFrame(transformed_sex_feature, columns = ["Sex_"+str(int(i)) for i in range(transformed_sex_feature.shape[1])])
df = pd.concat([df, df_sex_encoded], axis=1)
```

```
In [34]: df.head()
```

```
Out[34]:
```

	Sex	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	Rings	Age
0	M	0.455	0.365	0.095	0.5140	0.2245	0.1010	0.150	15	16.5
1	M	0.350	0.265	0.090	0.2255	0.0995	0.0485	0.070	7	8.5
2	F	0.530	0.420	0.135	0.6770	0.2565	0.1415	0.210	9	10.5
3	M	0.440	0.365	0.125	0.5160	0.2155	0.1140	0.155	10	11.5
4	I	0.330	0.255	0.080	0.2050	0.0895	0.0395	0.055	7	8.5

THANK YOU

```
In [ ]:
```