

```

1 import pandas as pd
2 import numpy as np
3 import statistics
4
5 data={'Name':['shreyash','hritik','sudhanva','ajay','vaibhav','subham'],
6       'Age':[12,17,22,18,24,30],
7       'Gender':['M','M','M','M','M','M'],
8       'Marks':[70,56,89,67,67,78],
9       'PhD':['Y','Y','N','Y','N','Y']}
10
11 df=pd.DataFrame(data)
12 df
13

```



	Name	Age	Gender	Marks	PhD
0	shreyash	12	M	70	Y
1	hritik	17	M	56	Y
2	sudhanva	22	M	89	N
3	ajay	18	M	67	Y
4	vaibhav	24	M	67	N
5	subham	30	F	78	Y

```

1
2 data2={'Name':['shreyash','hritik','sudhanva','ajay','vaibhav','subham'],
3        'Age':[12,17,22,18,np.NaN,30],
4        'Gender':['M','M','N/a','M','F','na'],
5        'Marks':[70,56,89,np.nan,67,78],
6        'PhD':['Y','Y','N',15,'N',np.nan]}
7 }
8 df2=pd.DataFrame(data2)
9 df2

```

	Name	Age	Gender	Marks	PhD
0	shreyash	12.0	M	70.0	Y
1	hritik	17.0	M	56.0	Y
2	sudhanva	22.0	N/a	89.0	N
3	ajay	18.0	M	NaN	15
4	vaibhav	NaN	F	67.0	N
5	subham	30.0	na	78.0	NaN

```

1 print (df2['Age'])
2 print(df2['Age'].isnull())

```

```

0    12.0
1    17.0
2    22.0
3    18.0
4     NaN
5    30.0
Name: Age, dtype: float64
0    False
1    False
2    False
3    False
4     True
5    False
Name: Age, dtype: bool

```

```

1 print(df2['Gender'])
2 print(df2['Gender'].isnull())

```

```

0     M
1     M
2    N/a
3     M
4     F
5    na
Name: Gender, dtype: object
0    False
1    False
2    False
3    False
4    False
5    False
Name: Gender, dtype: bool

```

```

1 #making list of missing values
2 #missing_values=['N/a','na']
3 #df2=pd.read_csv(df2,na_values=missing_values)
4 #print(df2['Gender'])
5 #print(df2['Gender'].isnull())
6

```

```

1 print(df2['PhD'])
2 print(df2['PhD'].isnull())

```

```

0     Y
1     Y
2     N
3    15
4     N
5    NaN
Name: PhD, dtype: object
0    False
1    False
2    False
3    False
4    False

```

```

5     True
Name: PhD, dtype: bool

1 #Detecting numbers
2 cnt=0
3 for row in df2['PhD']:
4     try:
5         int(row)
6         df2.loc[cnt, 'PhD']=np.nan
7     except ValueError:
8         pass
9     cnt+=1
10 print(df2['PhD'])
11 print(df2['PhD'].isnull())

```

```

0     Y
1     Y
2     N
3    NaN
4     N
5    NaN
Name: PhD, dtype: object
0    False
1    False
2    False
3     True
4    False
5     True
Name: PhD, dtype: bool

```

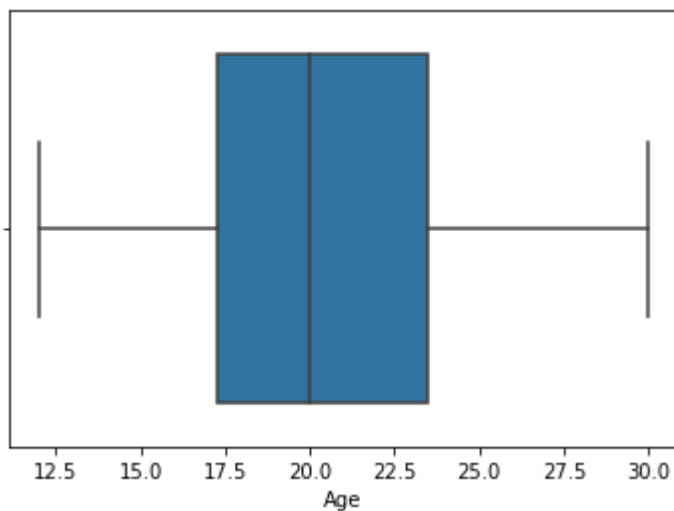
Double-click (or enter) to edit

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 sns.boxplot(x=df['Age'])

```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f842afa3bd0>



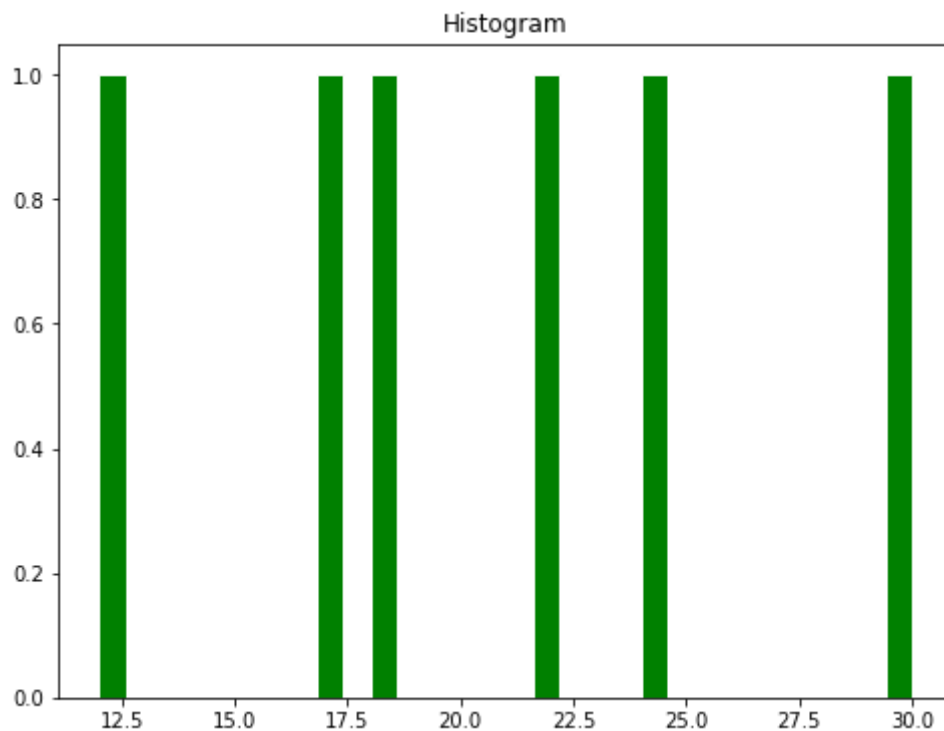
```

1 #position of outlier
2 print(np.where(df['Age']>20))

```

```
(array([2, 4, 5]),)
```

```
1 fig,x=plt.subplots(figsize=(8,6))
2 ax=plt.hist(df['Age'],bins=30,color='g',edgecolor='w')
3 plt.title('Histogram')
4 plt.show()
```



```
1 fig,ax=plt.subplots(figsize=(5,5))
2 ax.scatter(df['Age'],df['Marks'])
3
4 #x-axis label
5 ax.set_xlabel('Age')
6
7 #y- axis label
8 ax.set_ylabel('Marks')
9 plt.show()
```

```

on 1
1 df['Log_Age']=np.log(df['Age'])
2 df

```

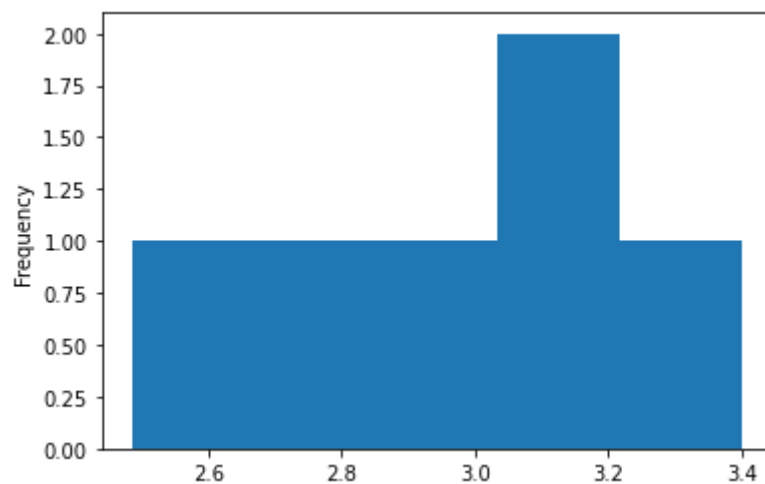
	Name	Age	Gender	Marks	PhD	Log_Age
0	shreyash	12	M	70	Y	2.484907
1	hritik	17	M	56	Y	2.833213
2	sudhanva	22	M	89	N	3.091042
3	ajay	18	M	67	Y	2.890372
4	vaibhav	24	M	67	N	3.178054
5	subham	30	F	78	Y	3.401197

```

1 df['Log_Age'].plot.hist(bins=5)

```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f842a87ff10>



```

1 df_scaled=df.copy()
2 col=['Age', 'Marks']
3 features=df_scaled[col]
4 from sklearn.preprocessing import MinMaxScaler
5 scaler=MinMaxScaler()
6 df_scaled[col]=scaler.fit_transform(features.values)
7 df_scaled

```

	Name	Age	Gender	Marko	BkD	Log Age
1						
1	hritik	0.277778	M	0.000000	Y	2.833213
2	sudhanva	0.555556	M	1.000000	N	3.091042
3	ajay	0.333333	M	0.333333	Y	2.890372
4	vaibhav	0.666667	M	0.333333	N	3.178054
5	subham	1.000000	F	0.666667	Y	3.401197

