

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
df=pd.read_csv('/content/water_quality.csv')
df
```

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity
0	NaN	204.890455	20791.318981	7.300212	368.516441	564.308654	10.379783	86.990970	2.96
1	3.716080	129.422921	18630.057858	6.635246	NaN	592.885359	15.180013	56.329076	4.50
2	8.099124	224.236259	19909.541732	9.275884	NaN	418.606213	16.868637	66.420093	3.05
3	8.316766	214.373394	22018.417441	8.059332	356.886136	363.266516	18.436524	100.341674	4.62
4	9.092223	181.101509	17978.986339	6.546600	310.135738	398.410813	11.558279	31.997993	4.07
...	...	...	...	...	...	...	...	...	...
3271	4.668102	193.681735	47580.991603	7.166639	359.948574	526.424171	13.894419	66.687695	2.96
3272	7.808856	193.553212	17329.802160	8.061362	NaN	392.449580	19.903225	NaN	4.50
3273	9.419510	175.762646	33155.578218	7.350233	NaN	432.044783	11.039070	69.845400	3.05
3274	5.126763	230.603758	11983.869376	6.303357	NaN	402.883113	11.168946	77.488213	4.62
3275	7.874671	195.102299	17404.177061	7.509306	NaN	327.459760	16.140368	78.698446	4.07

3276 rows × 10 columns

```
df.head()
```

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity
0	NaN	204.890455	20791.318981	7.300212	368.516441	564.308654	10.379783	86.990970	2.96
1	3.716080	129.422921	18630.057858	6.635246	NaN	592.885359	15.180013	56.329076	4.50
2	8.099124	224.236259	19909.541732	9.275884	NaN	418.606213	16.868637	66.420093	3.05
3	8.316766	214.373394	22018.417441	8.059332	356.886136	363.266516	18.436524	100.341674	4.62
4	9.092223	181.101509	17978.986339	6.546600	310.135738	398.410813	11.558279	31.997993	4.07

```
df.tail()
```

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity
3271	4.668102	193.681735	47580.991603	7.166639	359.948574	526.424171	13.894419	66.687695	2.96
3272	7.808856	193.553212	17329.802160	8.061362	NaN	392.449580	19.903225	NaN	4.50
3273	9.419510	175.762646	33155.578218	7.350233	NaN	432.044783	11.039070	69.845400	3.05
3274	5.126763	230.603758	11983.869376	6.303357	NaN	402.883113	11.168946	77.488213	4.62
3275	7.874671	195.102299	17404.177061	7.509306	NaN	327.459760	16.140368	78.698446	4.07

```
df.shape
```

(3276, 10)

```
df.dtypes
```

```
ph                float64
Hardness          float64
Solids            float64
Chloramines       float64
Sulfate           float64
Conductivity      float64
Organic_carbon    float64
Trihalomethanes   float64
Turbidity         float64
Potability        int64
dtype: object
```

```
df.isna().sum()
```

ph 491

```
Solids          0
Chloramines     0
Sulfate        781
Conductivity    0
Organic_carbon  0
Trihalomethanes 162
Turbidity       0
Potability      0
dtype: int64
```

df.describe()

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes
count	2785.000000	3276.000000	3276.000000	3276.000000	2495.000000	3276.000000	3276.000000	3114.000000
mean	7.080795	196.369496	22014.092526	7.122277	333.775777	426.205111	14.284970	66.396000
std	1.594320	32.879761	8768.570828	1.583085	41.416840	80.824064	3.308162	16.175000
min	0.000000	47.432000	320.942611	0.352000	129.000000	181.483754	2.200000	0.738000
25%	6.093092	176.850538	15666.690297	6.127421	307.699498	365.734414	12.065801	55.844000
50%	7.036752	196.967627	20927.833607	7.130299	333.073546	421.884968	14.218338	66.622000
75%	8.062066	216.667456	27332.762127	8.114887	359.950170	481.792304	16.557652	77.337000
max	14.000000	323.124000	61227.196008	13.127000	481.030642	753.342620	28.300000	124.000000

df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3276 entries, 0 to 3275
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   ph                    2785 non-null   float64
1   Hardness              3276 non-null   float64
2   Solids                3276 non-null   float64
3   Chloramines           3276 non-null   float64
4   Sulfate               2495 non-null   float64
5   Conductivity          3276 non-null   float64
6   Organic_carbon        3276 non-null   float64
7   Trihalomethanes       3114 non-null   float64
8   Turbidity             3276 non-null   float64
9   Potability            3276 non-null   int64
dtypes: float64(9), int64(1)
memory usage: 256.1 KB
```

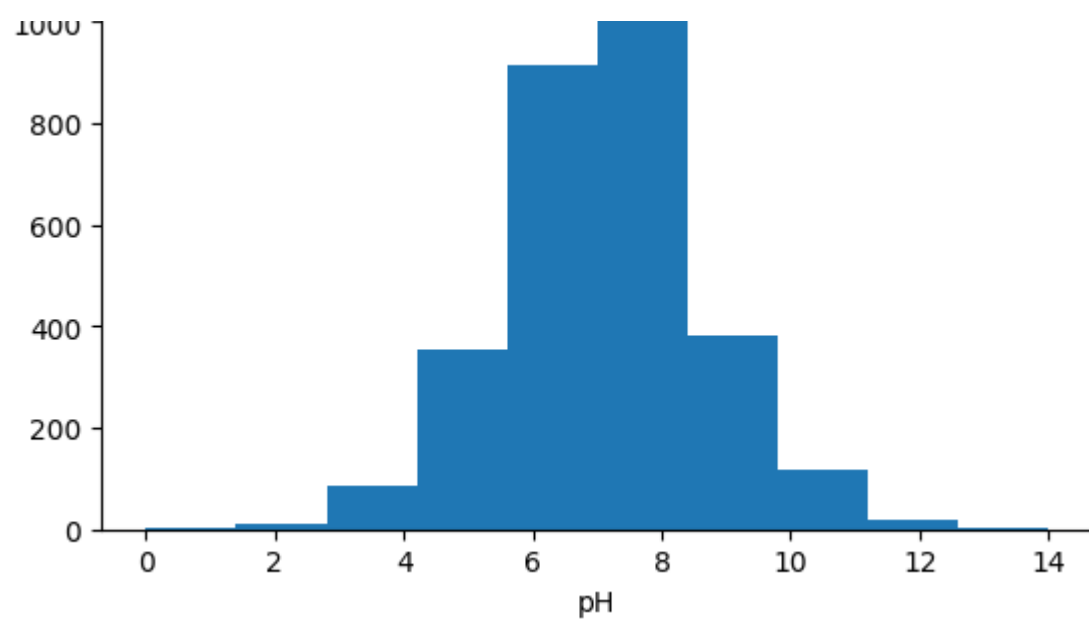
```
# Filling the missing values
df['ph'].fillna(df['ph'].mean(),inplace=True)
df['Sulfate'].fillna(df['Sulfate'].mean(),inplace=True)
df['Trihalomethanes'].fillna(df['Trihalomethanes'].mean(),inplace=True)
```

df.isna().sum()

```
ph          0
Hardness    0
Solids       0
Chloramines  0
Sulfate      0
Conductivity 0
Organic_carbon 0
Trihalomethanes 0
Turbidity    0
Potability   0
dtype: int64
```

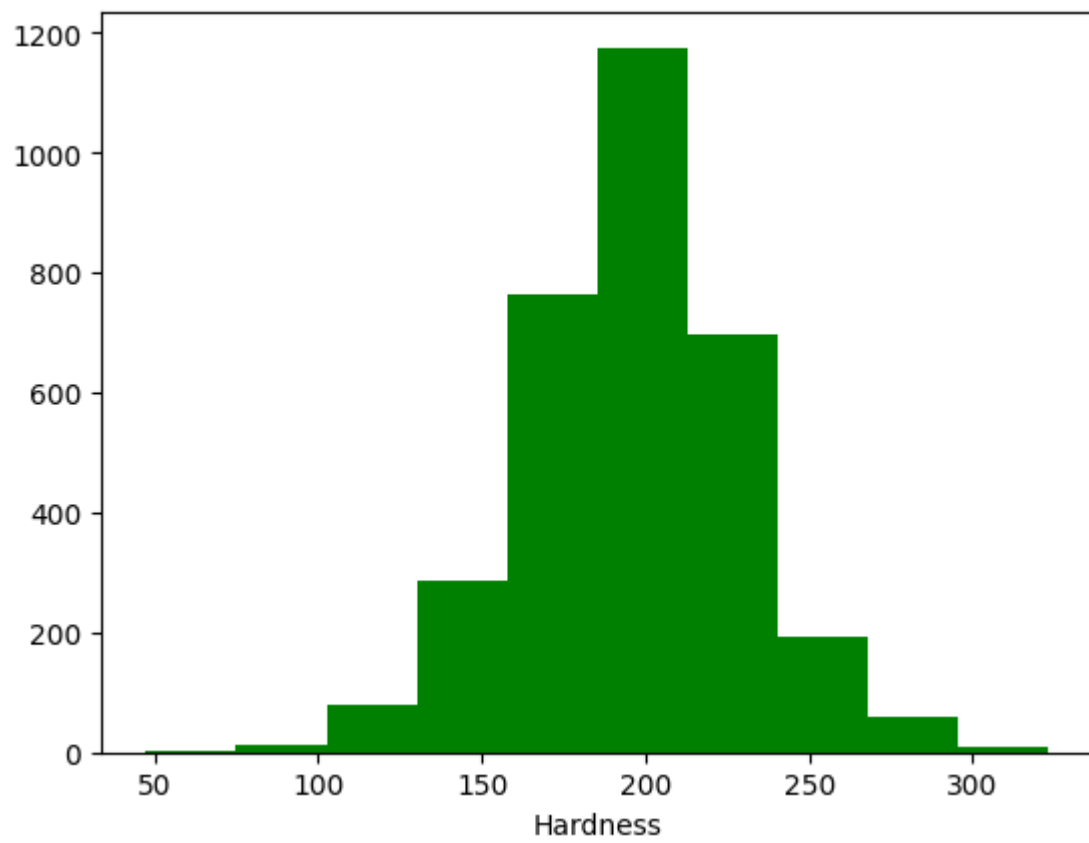
```
plt.hist(df['ph'])
plt.xlabel('pH')
```





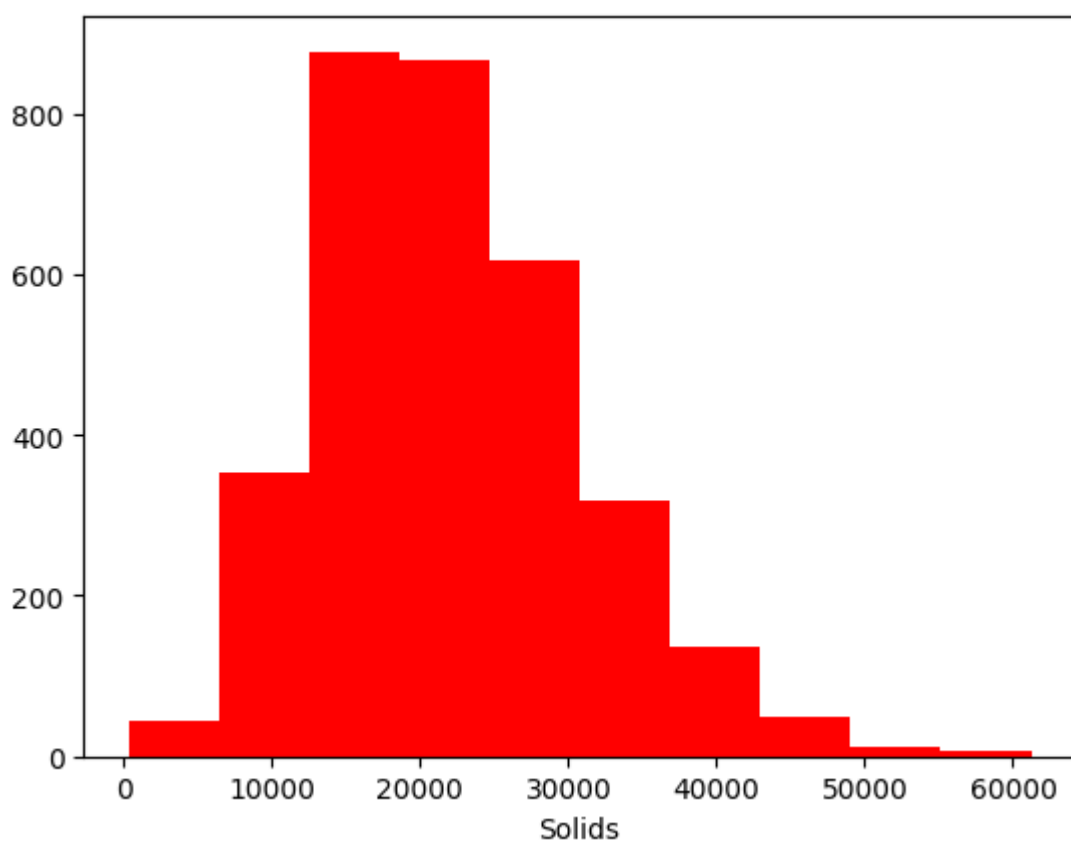
```
plt.hist(df['Hardness'],color='g')
plt.xlabel('Hardness')
```

Text(0.5, 0, 'Hardness')



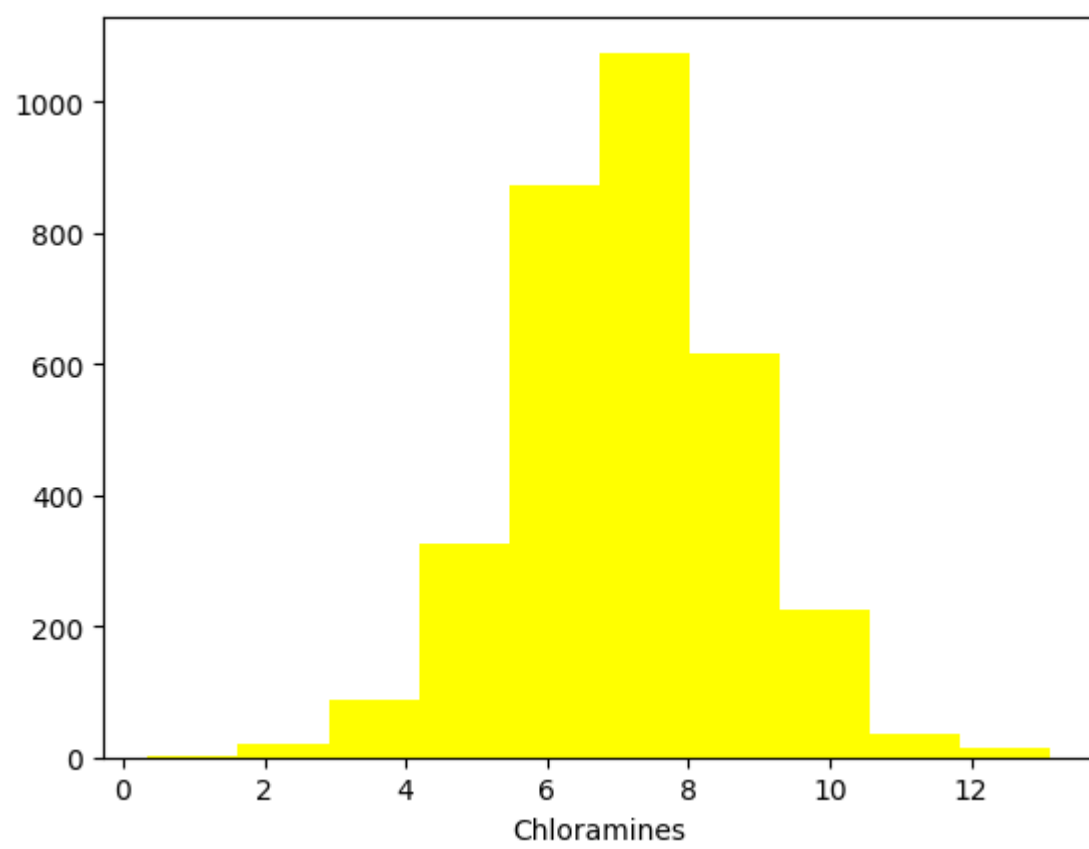
```
plt.hist(df['Solids'],color='r')
plt.xlabel('Solids')
```

Text(0.5, 0, 'Solids')



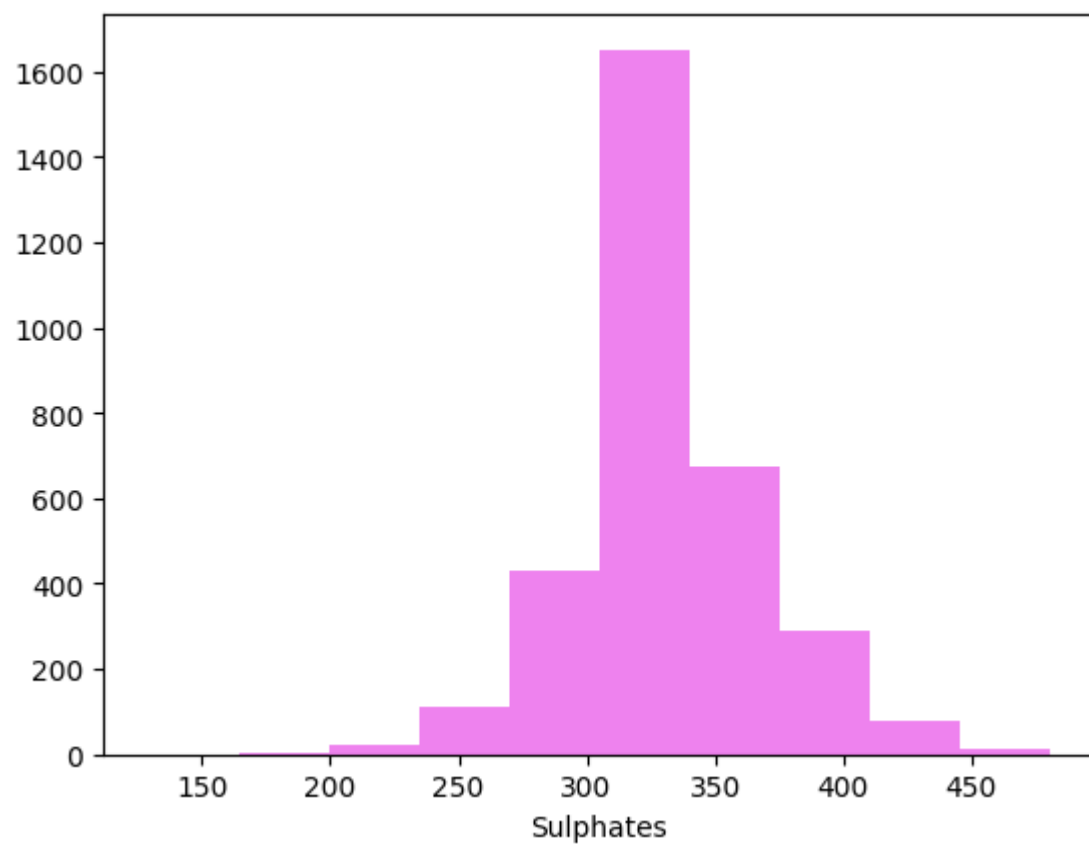
```
plt.hist(df['Chloramines'],color='yellow')
plt.xlabel('Chloramines')
```

Text(0.5, 0, 'Chloramines')



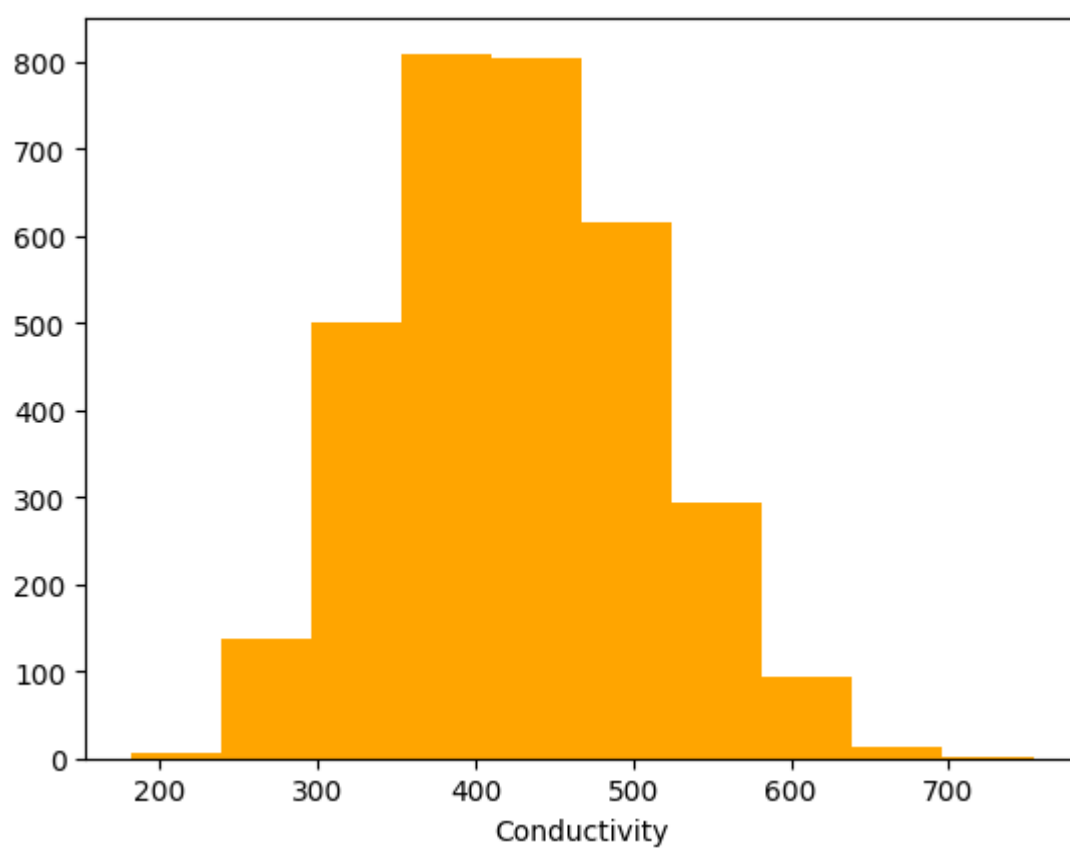
```
plt.hist(df['Sulfate'],color='violet')
plt.xlabel('Sulphates')
```

Text(0.5, 0, 'Sulphates')



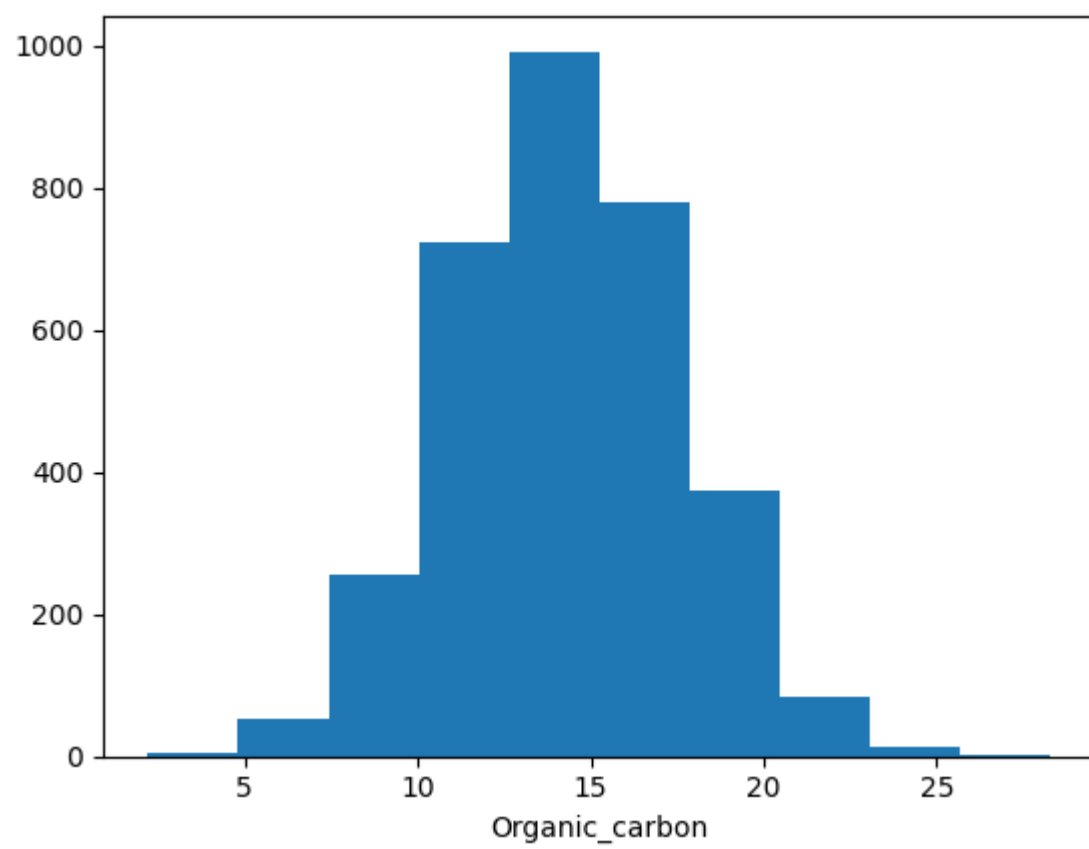
```
plt.hist(df['Conductivity'],color='orange')
plt.xlabel('Conductivity')
```

Text(0.5, 0, 'Conductivity')



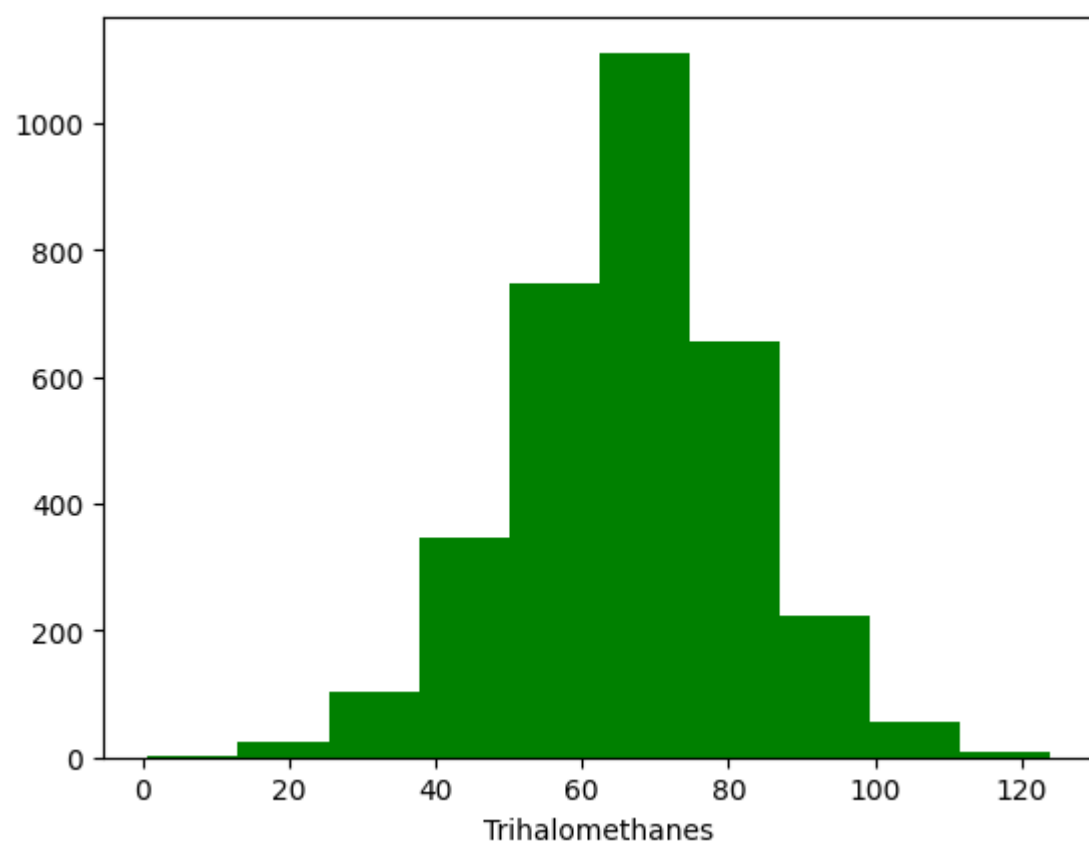
```
plt.hist(df['Organic_carbon'])  
plt.xlabel('Organic_carbon')
```

Text(0.5, 0, 'Organic\_carbon')



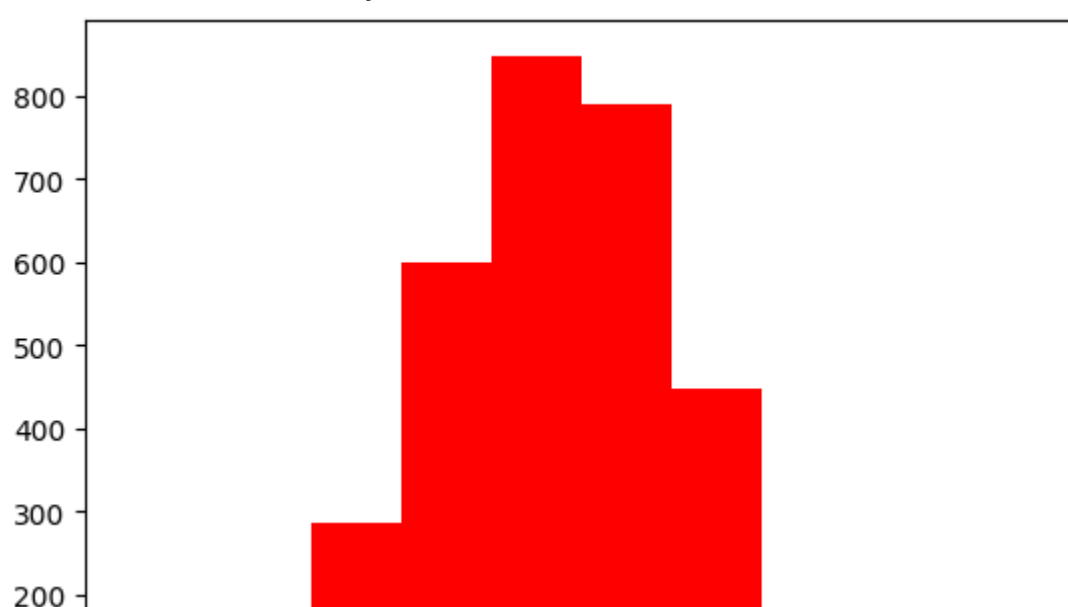
```
plt.hist(df['Trihalomethanes'],color='g')  
plt.xlabel('Trihalomethanes')
```

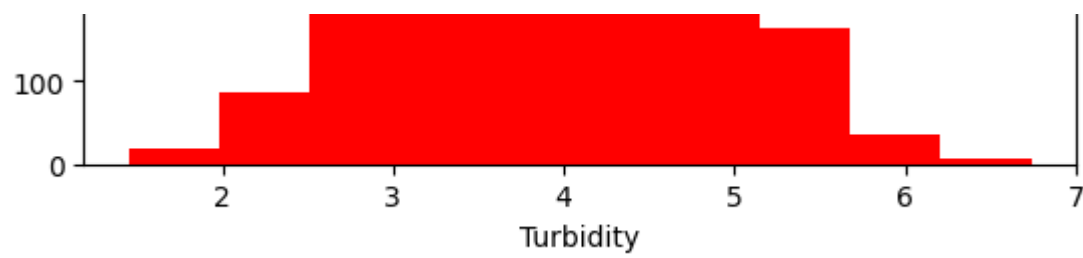
Text(0.5, 0, 'Trihalomethanes')



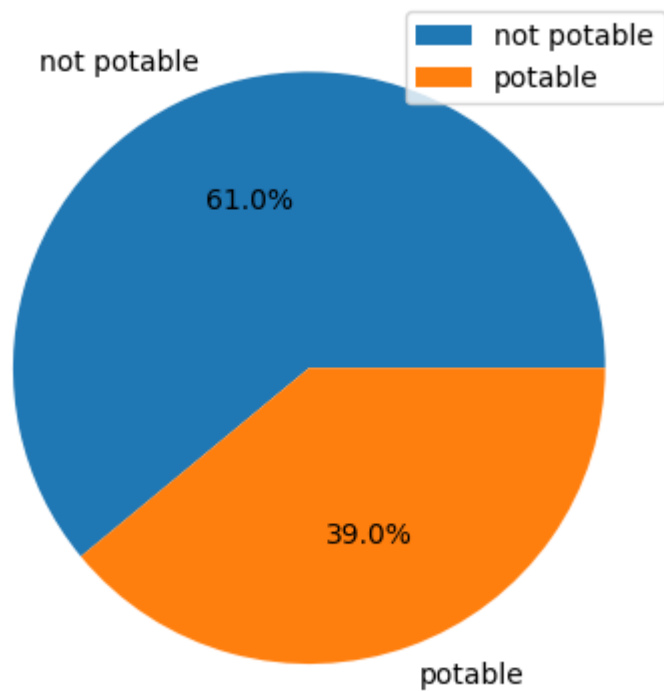
```
plt.hist(df['Turbidity'],color='r')  
plt.xlabel('Turbidity')
```

Text(0.5, 0, 'Turbidity')





```
df1=df['Potability'].value_counts()
potability=['not potable','potable']
plt.pie(df1,labels=potability,autopct='%1.1f%%')
plt.legend()
plt.show()
df1
```



```
0    1998
1    1278
Name: Potability, dtype: int64
```

## Separating input and output samples

```
x=df.iloc[:, :-1].values
y=df.iloc[:, -1].values
x,y
```

```
(array([[7.08079450e+00, 2.04890455e+02, 2.07913190e+04, ...,
        1.03797831e+01, 8.69909705e+01, 2.96313538e+00],
       [3.71608008e+00, 1.29422921e+02, 1.86300579e+04, ...,
        1.51800131e+01, 5.63290763e+01, 4.50065627e+00],
       [8.09912419e+00, 2.24236259e+02, 1.99095417e+04, ...,
        1.68686369e+01, 6.64200925e+01, 3.05593375e+00],
       ...,
       [9.41951032e+00, 1.75762646e+02, 3.31555782e+04, ...,
        1.10390697e+01, 6.98454003e+01, 3.29887550e+00],
       [5.12676292e+00, 2.30603758e+02, 1.19838694e+04, ...,
        1.11689462e+01, 7.74882131e+01, 4.70865847e+00],
       [7.87467136e+00, 1.95102299e+02, 1.74041771e+04, ...,
        1.61403676e+01, 7.86984463e+01, 2.30914906e+00]]),
 array([0, 0, 0, ..., 1, 1, 1]))
```

## Training and testing data

```
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.30,random_state=42)
```

```
x_train
```

```
array([[7.08079450e+00, 1.88445469e+02, 2.87916144e+04, ...,
        1.05756901e+01, 6.32353650e+01, 3.22837922e+00],
       [7.20343885e+00, 1.68445358e+02, 2.28264847e+04, ...,
        1.64106541e+01, 6.45059226e+01, 6.38916101e+00],
       [7.08079450e+00, 2.42827588e+02, 2.92980743e+04, ...,
        5.42664993e+00, 6.63962929e+01, 3.52258617e+00],
```

x\_test

y\_train

y\_test

## Normalisation

```
from sklearn.preprocessing import MinMaxScaler
scaler=MinMaxScaler()
```





## Performance evaluation

```
from sklearn.metrics import classification_report
print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.68	0.95	0.79	617
1	0.75	0.23	0.35	366
accuracy			0.69	983
macro avg	0.71	0.59	0.57	983
weighted avg	0.71	0.69	0.63	983