

DATA PREPROCESSING IN PYTHON

Steps of preprocessing of data

Step-1: IMPORT NECESSARY LIBRARY

Step-2: READ DATASET

Step-3: SANITY CHECK OF DATA

Step-4: EXPLORATORY DATA ANALYSIS (EDA)

Step-5: MISSING VALUE TREATMENTS

Step-6: OUTLIERS TREATMENTS

Step-7: DUPLICATE VALUES AND GARBAGE VALUES TREATMENT

STEP-1 IMPORT NECESSARY LIBRARY

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

STEP-2 READ DATASET

```
df=pd.read_csv('D:\\Air Quality Data\\Air_Cleaned_Processed.csv')
df.head() #It will show top 5 data
```

	City	Date	PM2.5	PM10	N02
0	Allahabad	2013-01-01	72.190988	28.038859	74.466703
1	Nagpur	2013-01-02	156.162525	109.705592	79.357025
2	Vadodara	2013-01-03	37.035248	320.996249	19.526083
3	Nashik	2013-01-04	222.000056	217.004109	85.572290
4	Srinagar	2013-01-05	259.690694	47.634851	135.613423

32.026520

	C0	03	Temperature	Humidity	AQI \
0	2.737498	22.832620	23.538722	91.462662	139.203717
1	8.298729	87.444488	31.233321	55.057146	327.985193
2	0.268302	14.381349	39.919461	86.670590	270.996249
3	1.372084	25.333474	39.880965	76.203614	378.511671
4	6.027527	95.918899	24.098538	75.658559	404.455336

	AQI Category	PM2.5/PM10 Ratio
0	Unhealthy for Sensitive Groups	1.500000
1	Hazardous	1.423469
2	Very Unhealthy	0.115376
3	Hazardous	1.023022
4	Hazardous	1.500000

df.tail() *#It will show bottom 5 data*

	City	Date	PM2.5	PM10	N02
S02 \					
3995	Kota	2023-12-10	216.579951	50.677063	74.476868
99.373330					
3996	Ranchi	2023-12-11	11.870857	239.002070	39.937385
5.815086					
3997	Howrah	2023-12-12	189.599794	312.727515	97.029882
18.905951					
3998	Hyderabad	2023-12-13	173.093545	53.921186	130.987842
97.598887					
3999	Jaipur	2023-12-14	275.059510	273.342239	118.910881
82.707042					

	C0	03	Temperature	Humidity	AQI	AQI
Category \						
3995	0.139059	57.206447	40.822590	77.312612	500.000000	
Hazardous						
3996	6.515407	79.966366	15.951747	80.208933	192.692651	
Unhealthy						
3997	7.323106	99.748538	10.104865	80.187652	353.646353	
Hazardous						
3998	2.455181	44.449088	20.986299	50.480373	340.978767	
Hazardous						
3999	7.439409	79.277567	31.209837	28.386461	410.565829	
Hazardous						

	PM2.5/PM10 Ratio
3995	1.500000
3996	0.049668
3997	0.606278
3998	1.500000
3999	1.006282

STEP-3 SANITY CHECK OF DATA

#SHAPE

```
df.shape
```

```
(4000, 13)
```

#Info

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 4000 entries, 0 to 3999
```

```
Data columns (total 13 columns):
```

#	Column	Non-Null Count	Dtype
0	City	4000 non-null	object
1	Date	4000 non-null	object
2	PM2.5	4000 non-null	float64
3	PM10	4000 non-null	float64
4	N02	4000 non-null	float64
5	S02	4000 non-null	float64
6	C0	4000 non-null	float64
7	O3	4000 non-null	float64
8	Temperature	4000 non-null	float64
9	Humidity	4000 non-null	float64
10	AQI	4000 non-null	float64
11	AQI Category	4000 non-null	object
12	PM2.5/PM10 Ratio	4000 non-null	float64

```
dtypes: float64(10), object(3)
```

```
memory usage: 406.4+ KB
```

#Finding Missing Value

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

City	0
Date	0
PM2.5	0
PM10	0
N02	0
S02	0
C0	0
O3	0
Temperature	0
Humidity	0
AQI	0
AQI Category	0
PM2.5/PM10 Ratio	0

```
dtype: int64
```

#Finding percentage of missing values

```
df.isnull().sum()/df.shape[0]*100
```

```
City          0.0
Date          0.0
PM2.5        0.0
PM10         0.0
NO2          0.0
SO2          0.0
CO           0.0
O3           0.0
Temperature  0.0
Humidity     0.0
AQI          0.0
AQI Category 0.0
PM2.5/PM10 Ratio 0.0
dtype: float64
```

#Finding Duplicates

```
df.duplicated().sum()
```

```
0
```

#FIXING CITY NAME ISSUE

```
df["City"] = df["City"].replace({"LucknowUntitled document":  
"Lucknow"})
```

Standardize AQI Category Formatting

```
df["AQI Category"] = df["AQI Category"].str.strip().str.title()  
#.str.strip() removes leading/trailing spaces.  
#.str.title() ensures consistent casing (e.g., "very unhealthy" →  
"Very Unhealthy").
```

#Date formatting

```
df["Date"] = pd.to_datetime(df["Date"], errors="coerce")  
#Converts Date to a datetime format.  
#If there are invalid dates, they become NaT (Not a Time), making them  
easy to identify
```

#Identifying Garbage Value

```
#Garbage values are always present in the form of object data type  
for i in df.select_dtypes(include="object").columns:  
    print(df[i].value_counts())  
    print("****"*10)
```

```
City
Chennai      116
Bangalore    114
Allahabad    111
Jabalpur     110
Kanpur       109
Chandigarh   109
Dhanbad      109
Raipur       109
```

Aurangabad	108
Solapur	108
Nashik	105
Bhopal	104
Mumbai	104
Varanasi	104
Surat	103
Jaipur	103
Vadodara	103
Vijayawada	103
Ludhiana	102
Patna	102
Hyderabad	100
Nagpur	100
Lucknow	99
Amritsar	99
Pune	99
Ahmedabad	98
Srinagar	98
Delhi	98
Agra	97
Ranchi	94
Kota	94
Howrah	93
Visakhapatnam	91
Jodhpur	90
Navi Mumbai	88
Kolkata	87
Gwalior	87
Madurai	86
Coimbatore	83
Guwahati	83

Name: count, dtype: int64

AQI Category	
Hazardous	2827
Very Unhealthy	581
Unhealthy	431
Unhealthy For Sensitive Groups	153
Moderate	8

Name: count, dtype: int64

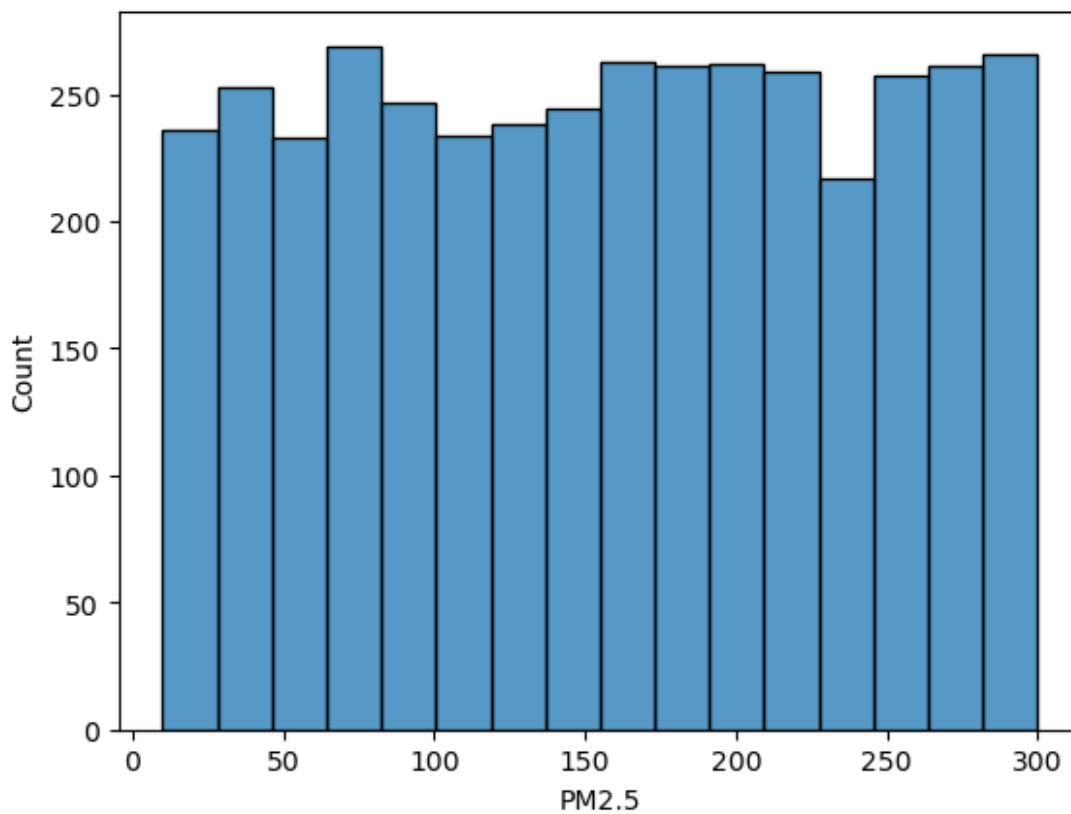
Step-4 EXPLORATORY DATA ANALYSIS (EDA)

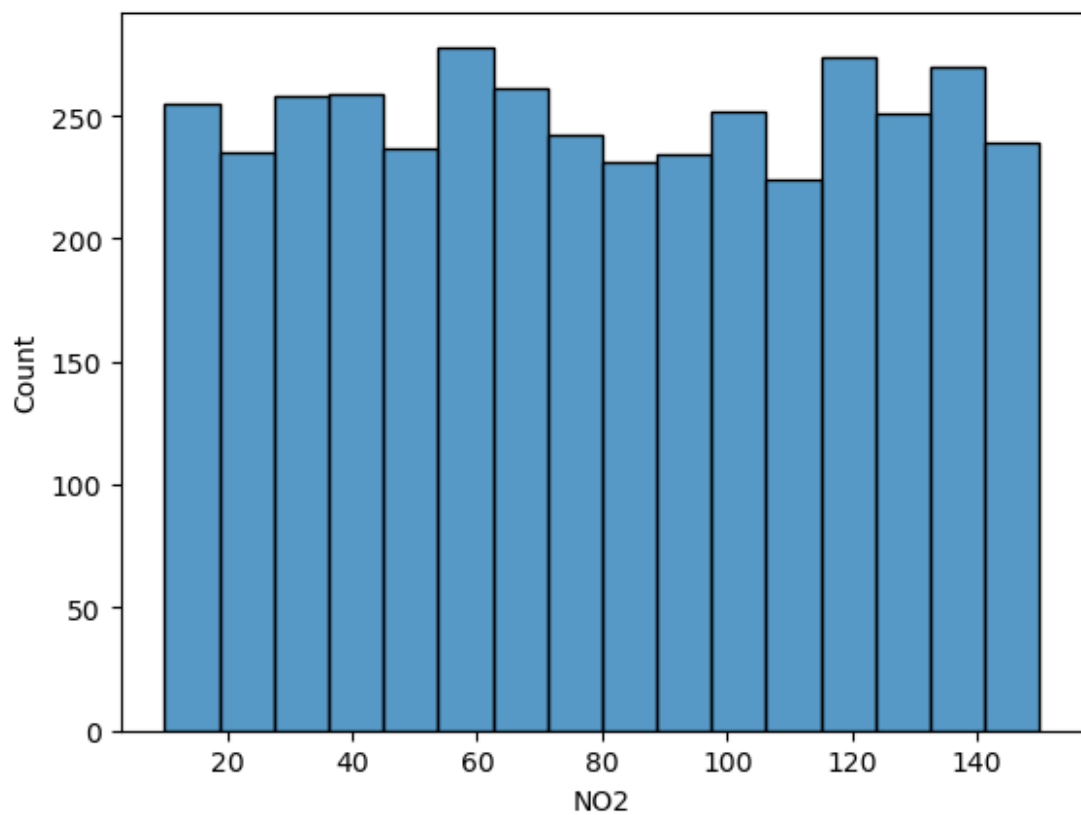
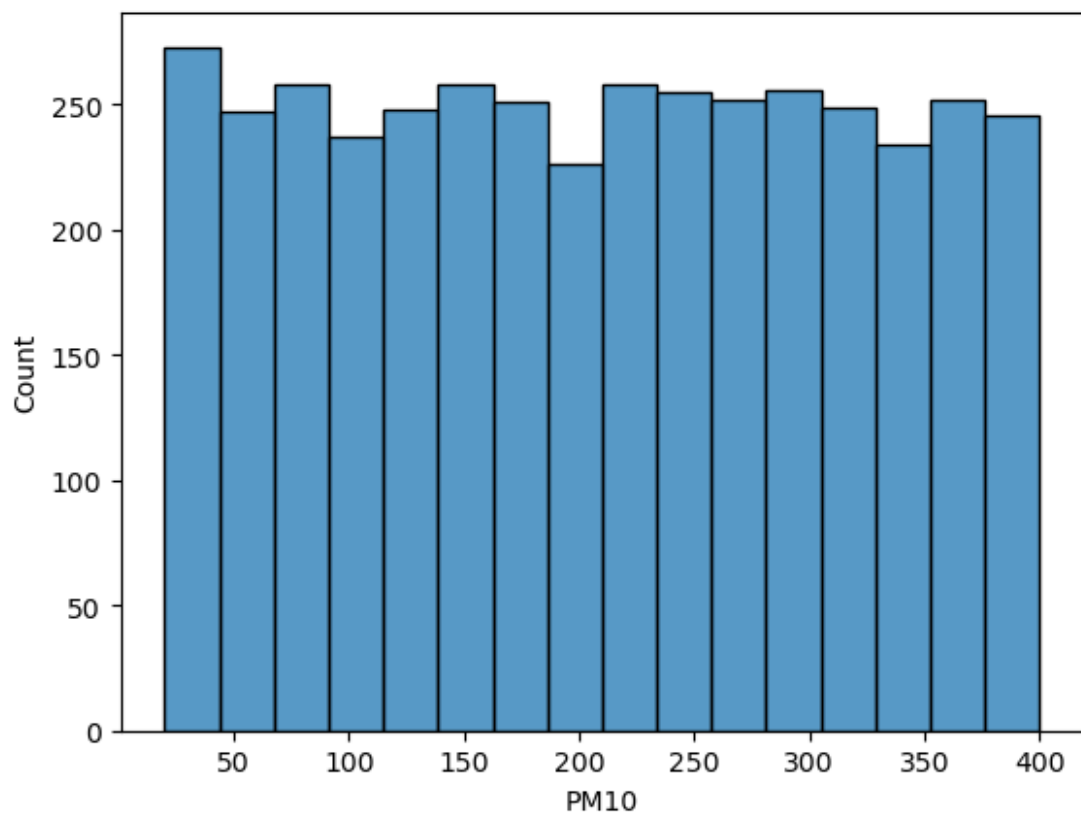
```
#Descriptive Statistics
# df.describe()
df.describe(include="object")
```

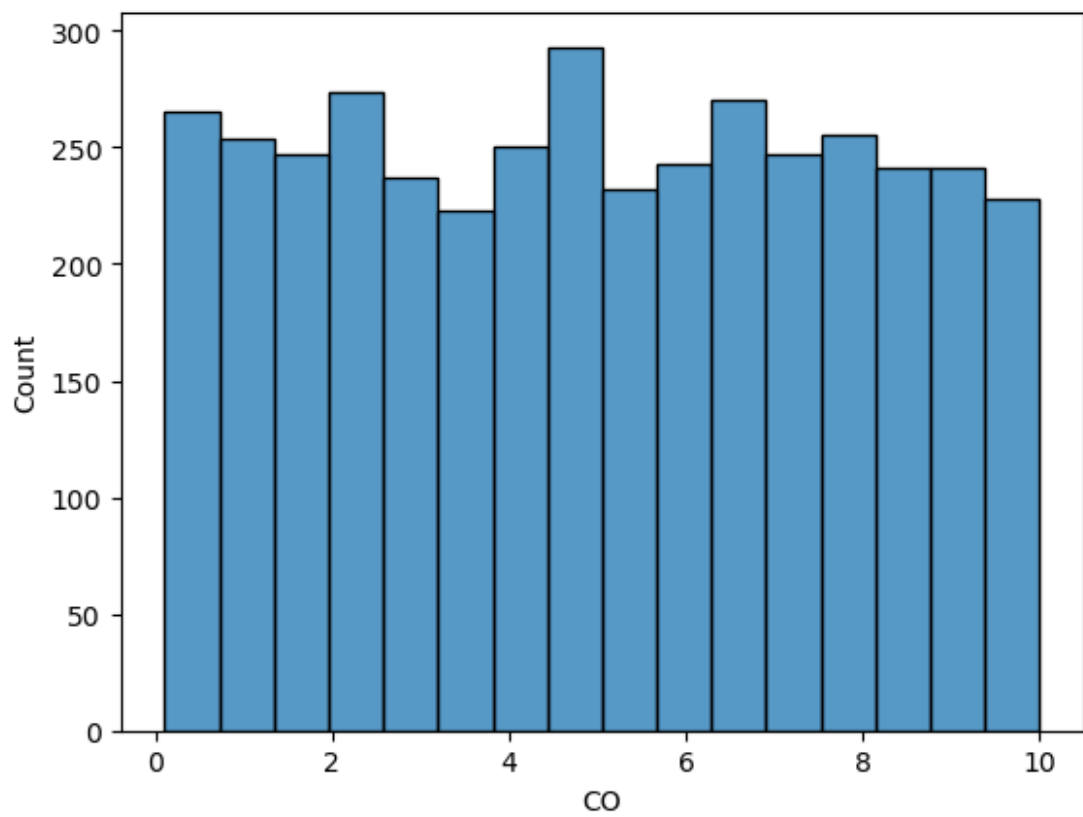
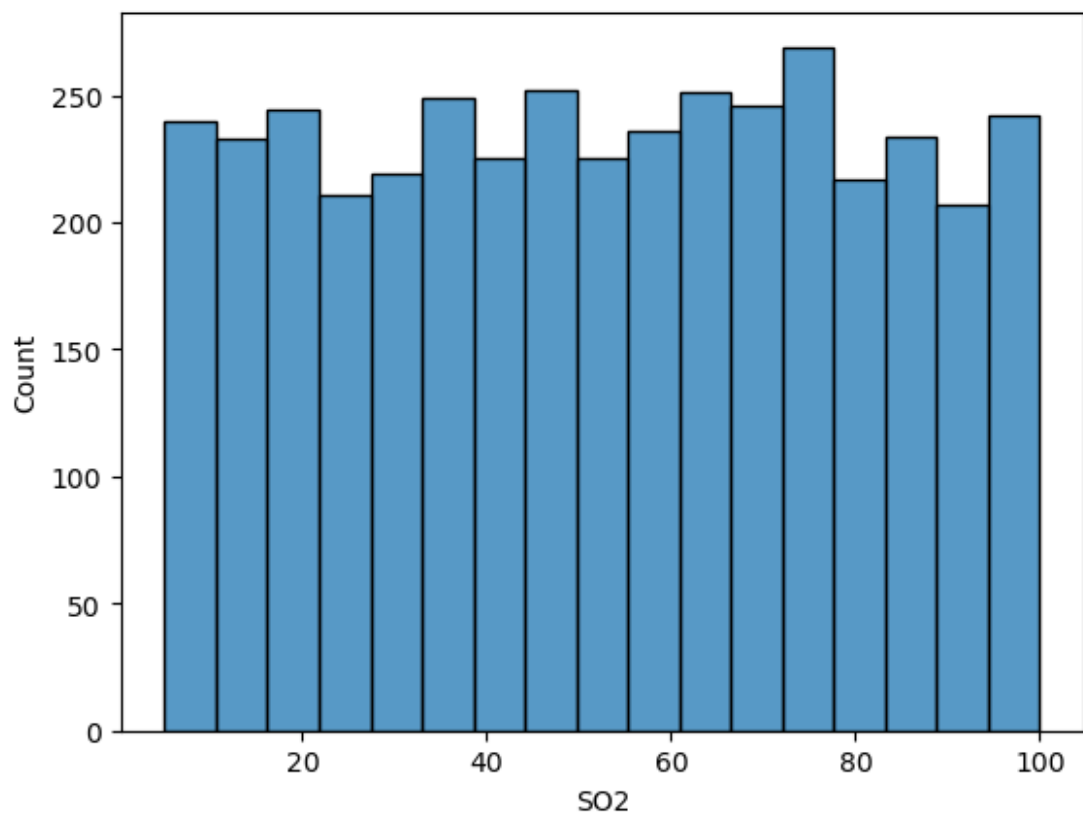
	City	AQI	Category
count	4000		4000
unique	40		5
top	Chennai		Hazardous
freq	116		2827

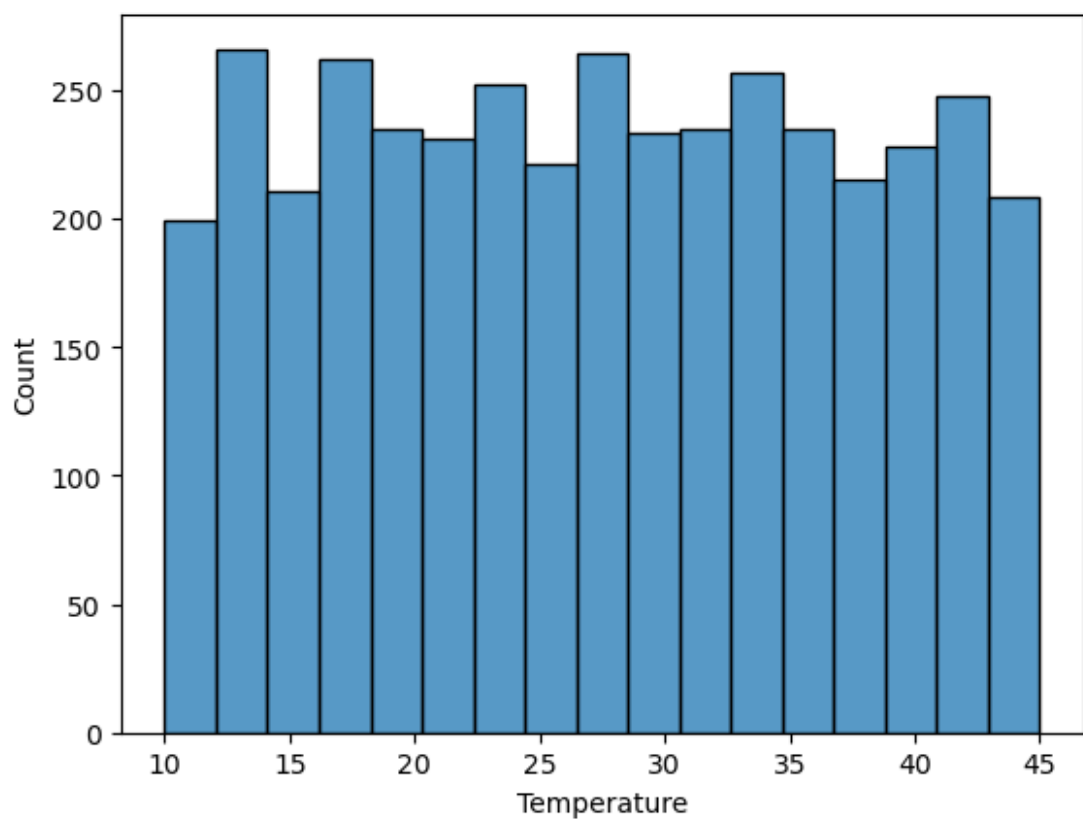
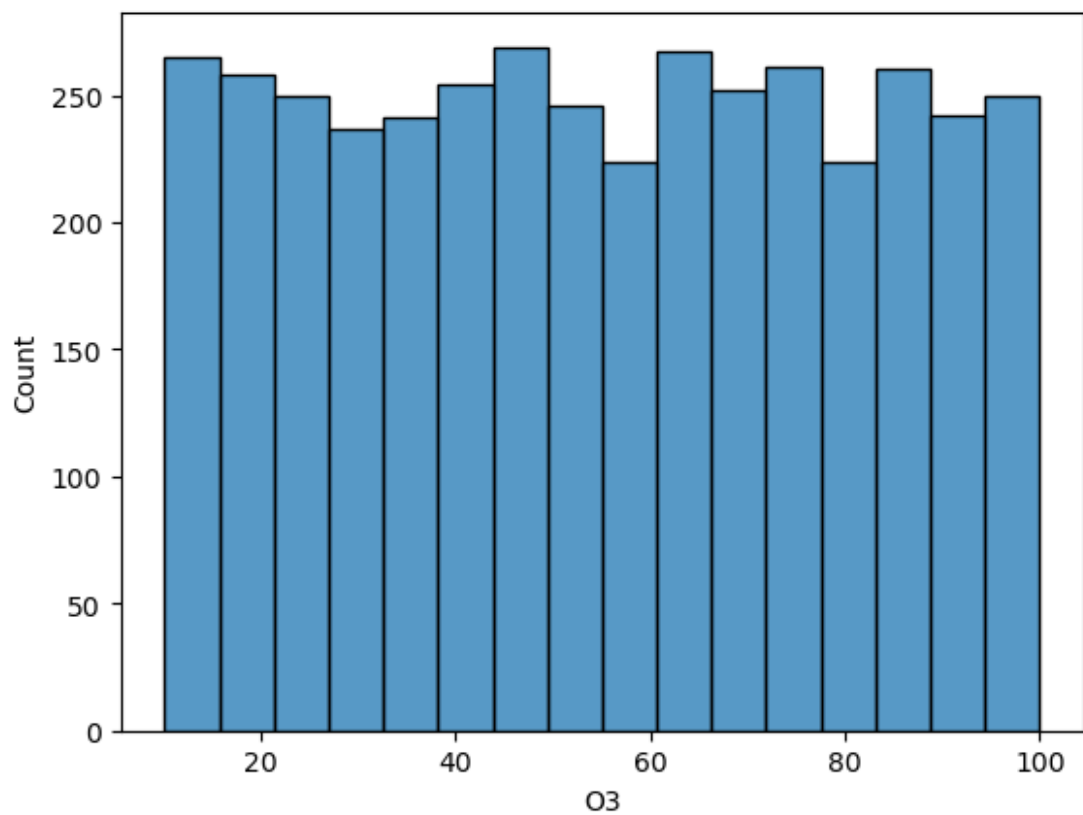
#Histogram to understand the distribution

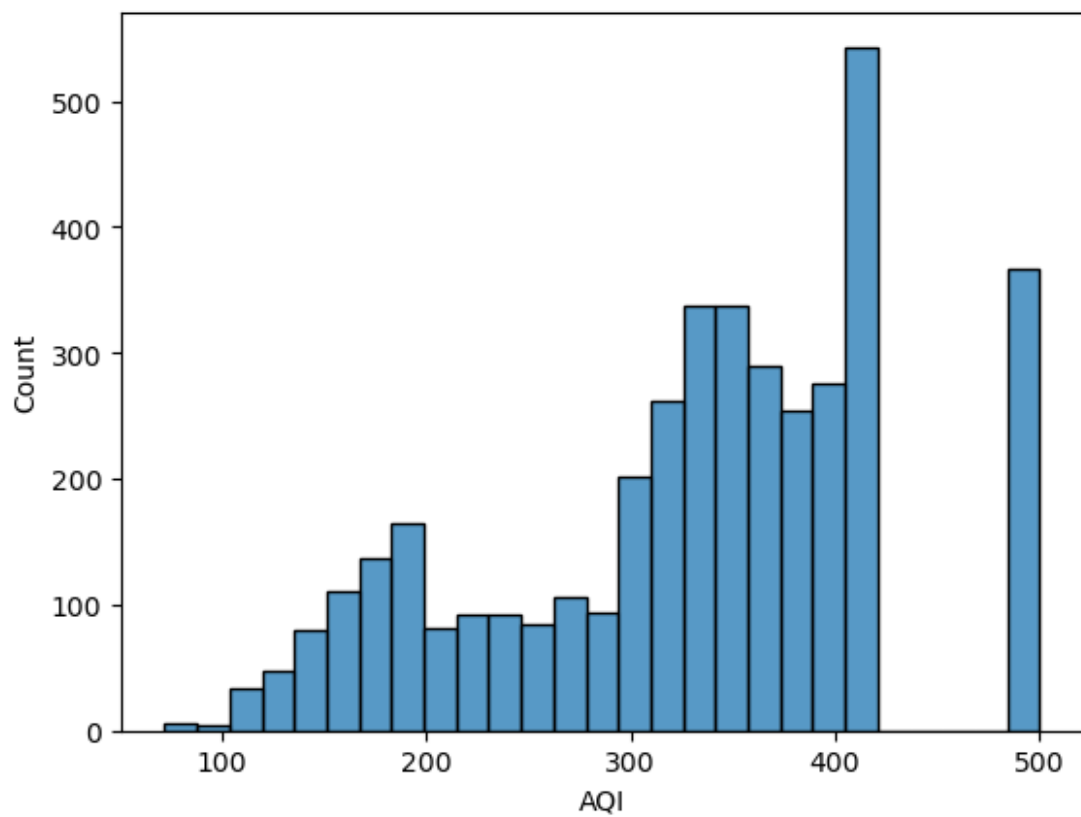
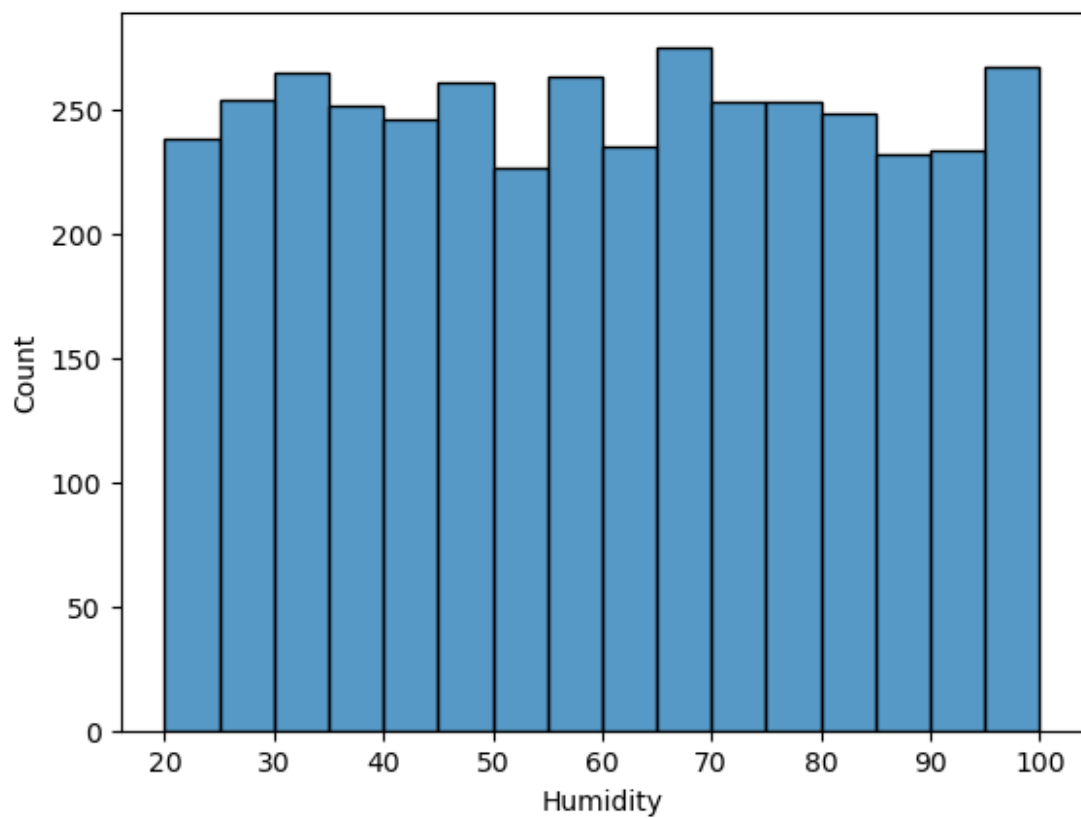
```
import warnings
warnings.filterwarnings("ignore")
for i in df.select_dtypes(include="number").columns:
    sns.histplot(data=df, x=i)
    plt.show()
```

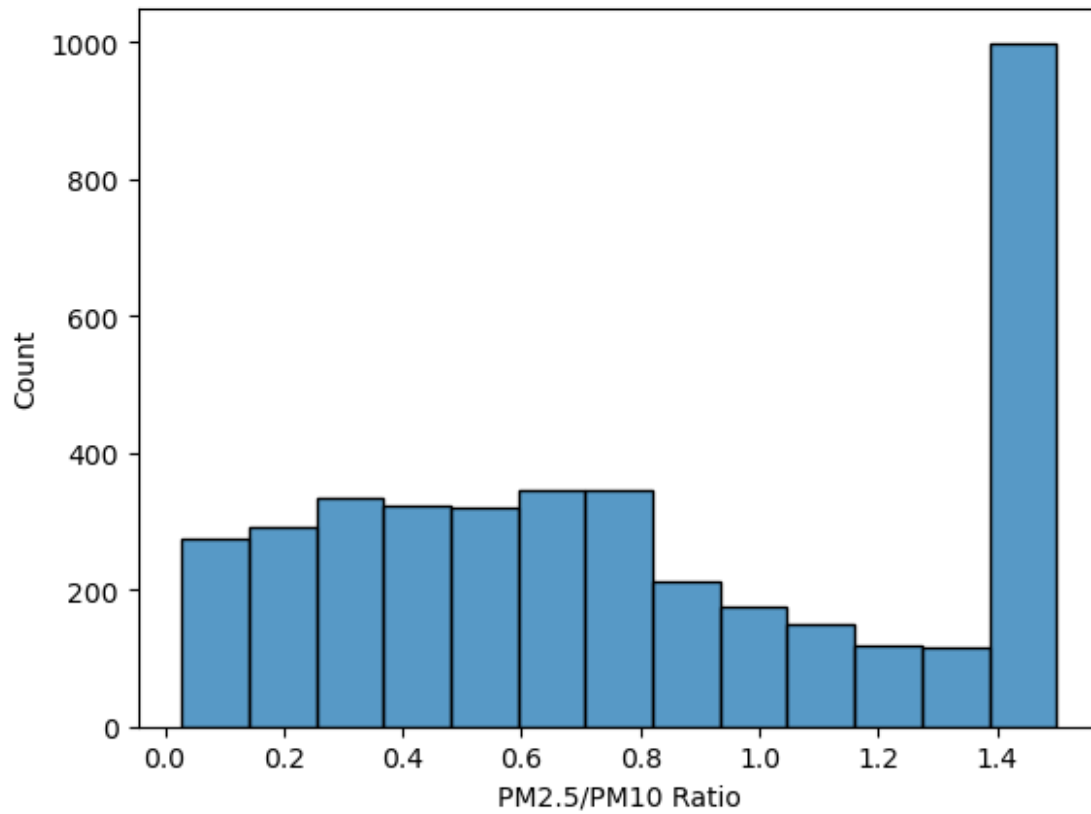




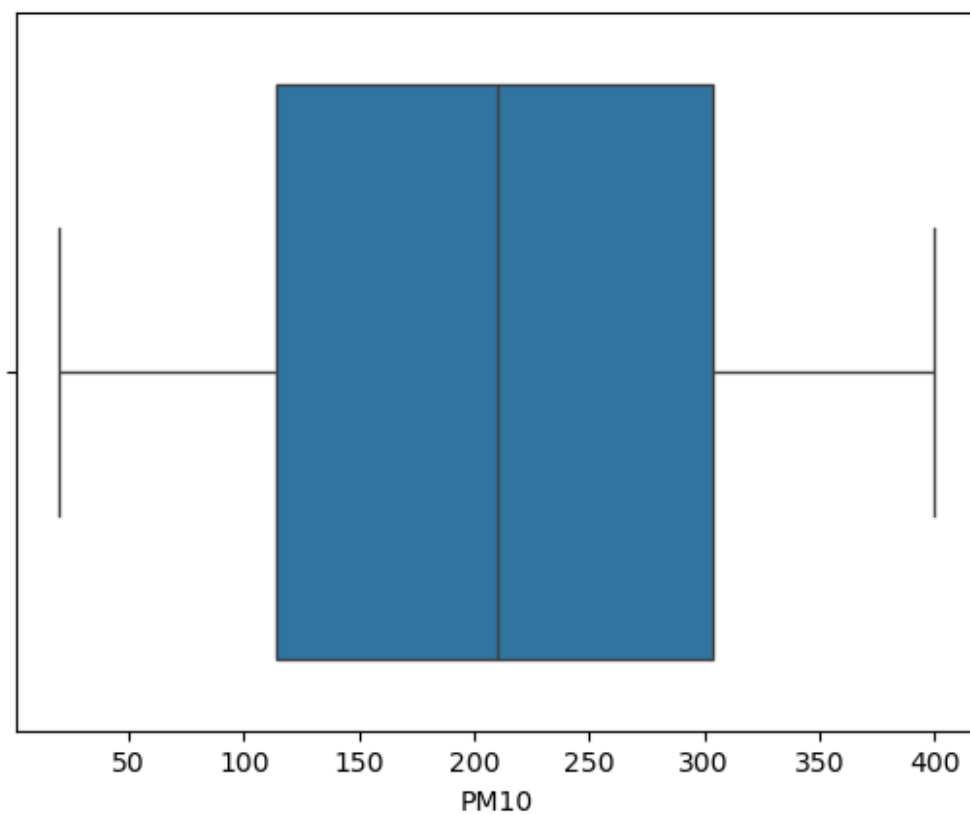
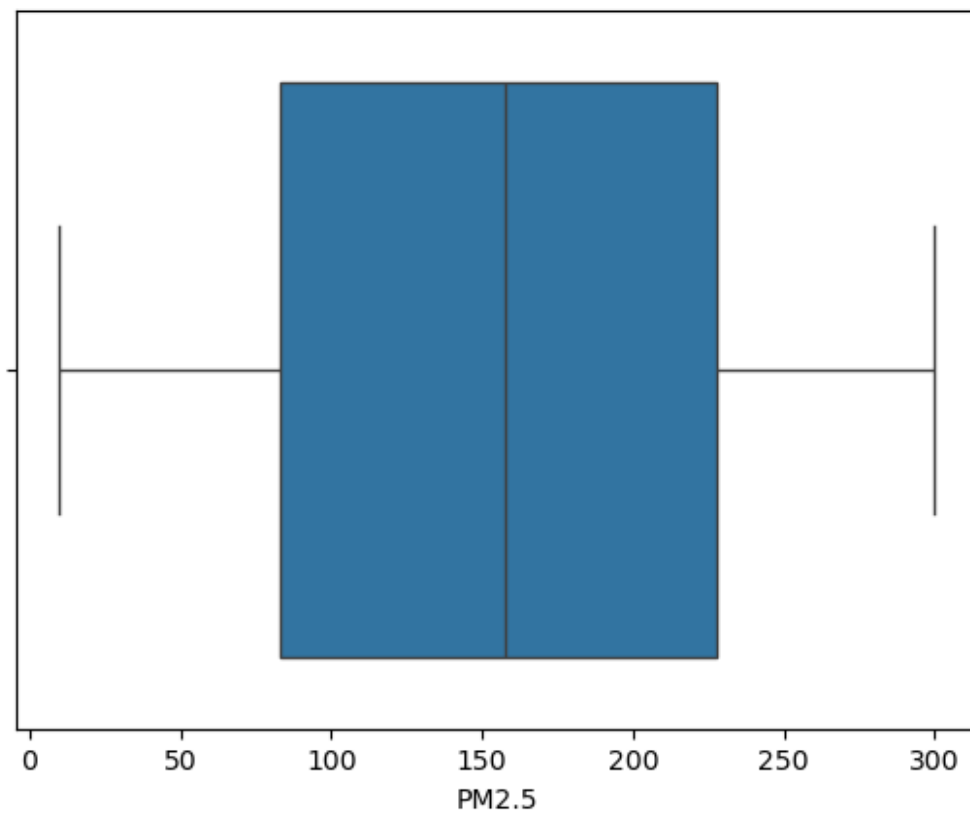


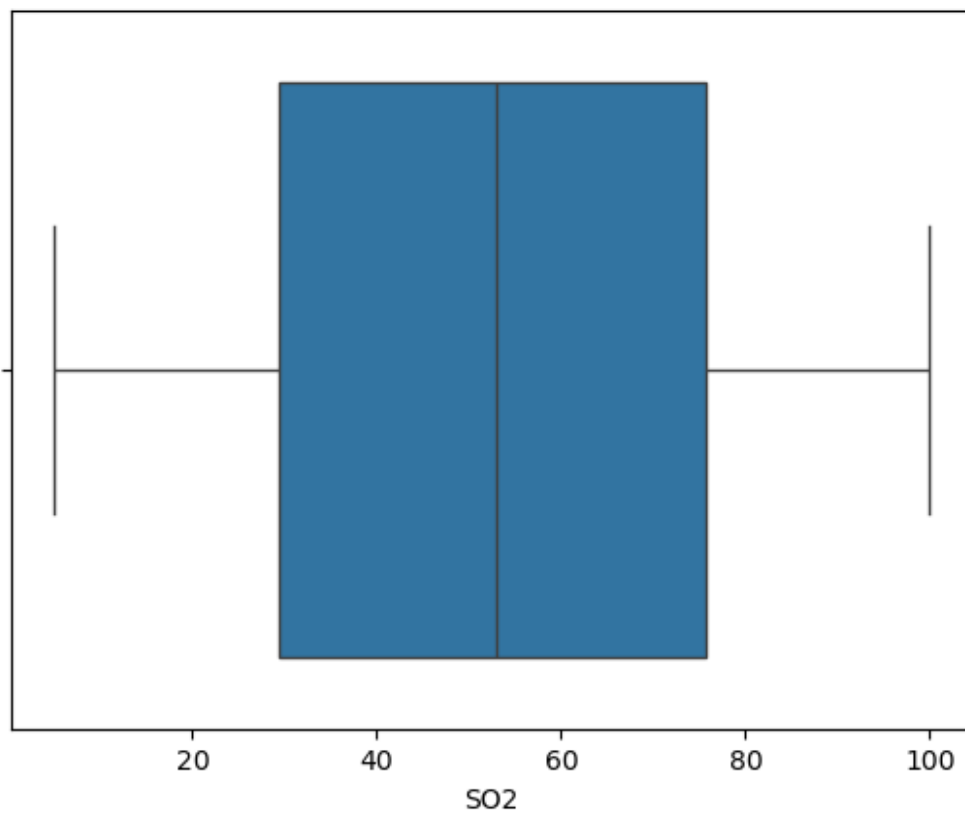
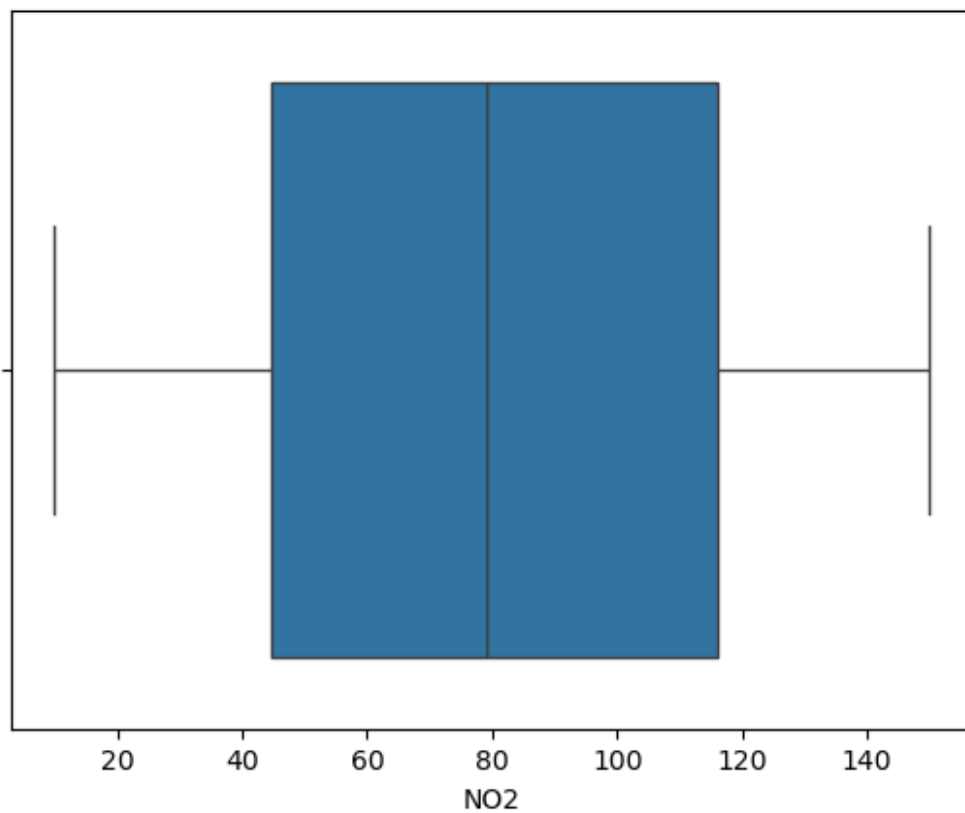


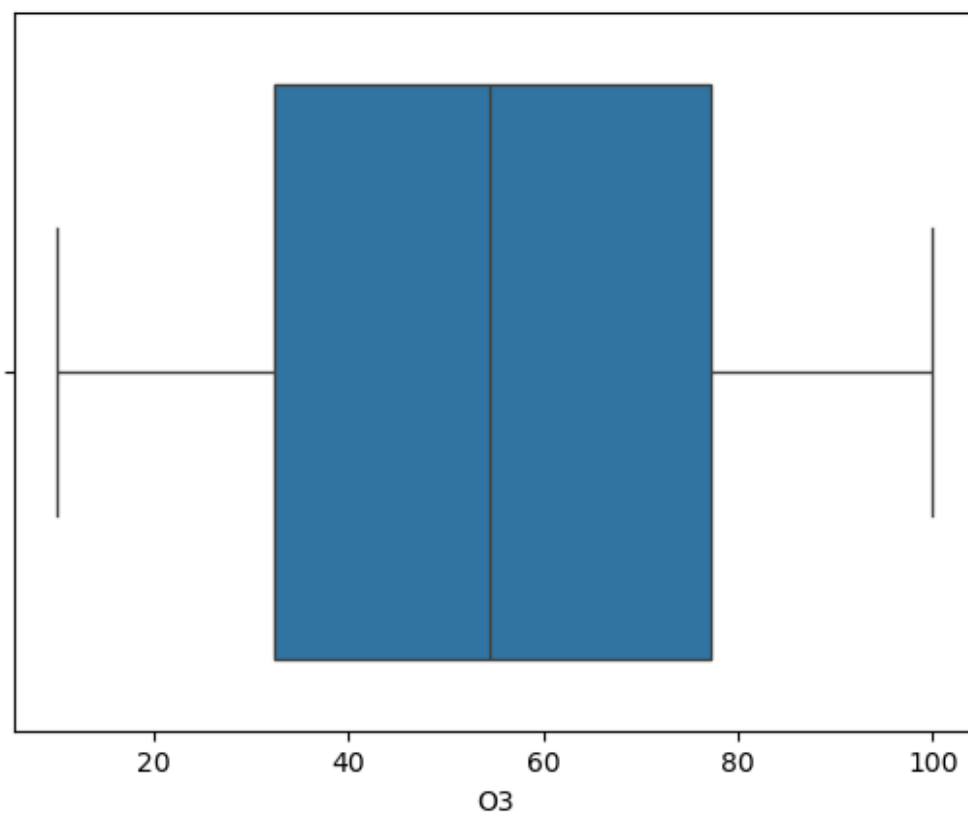
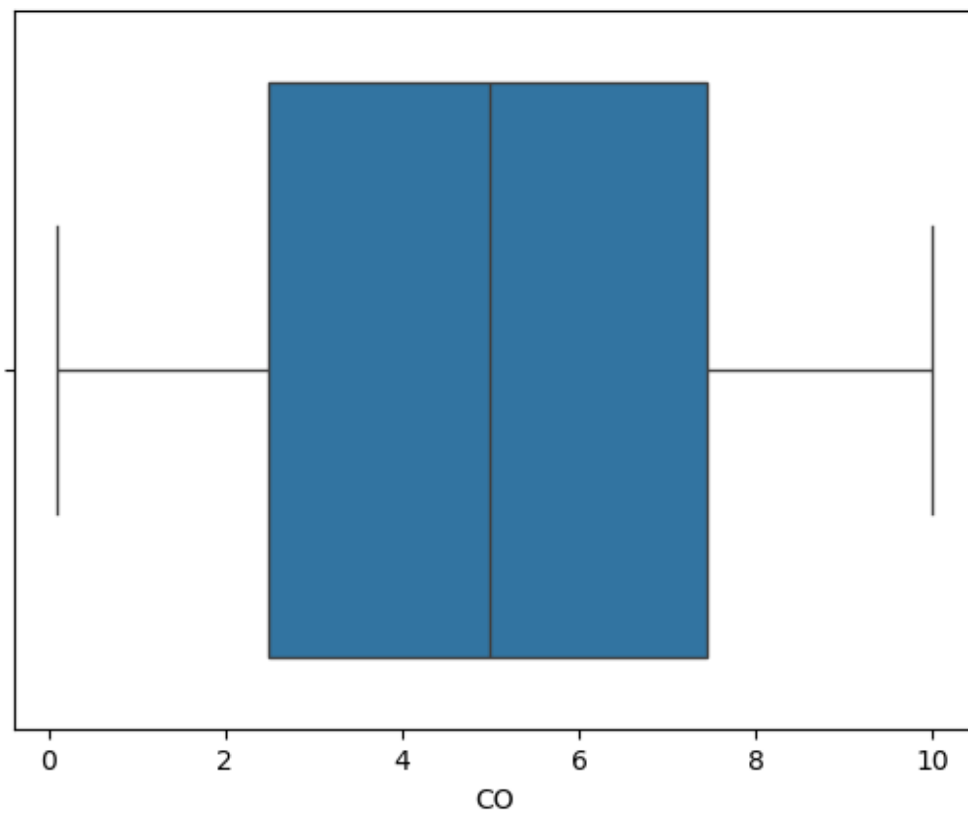


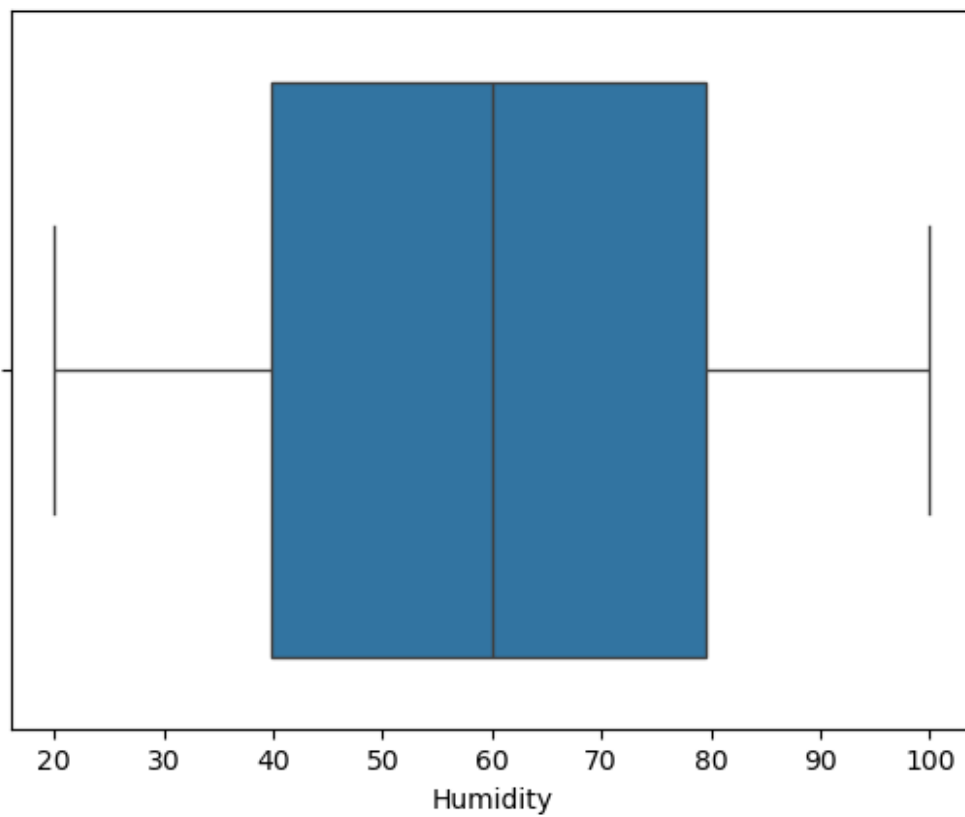
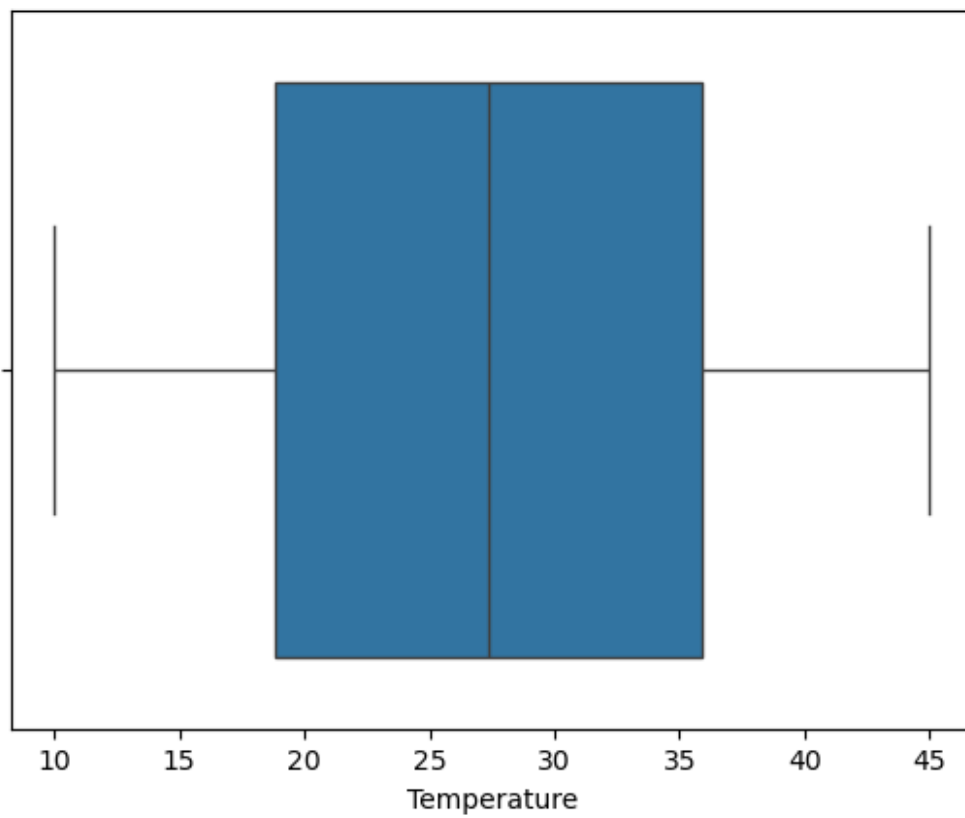


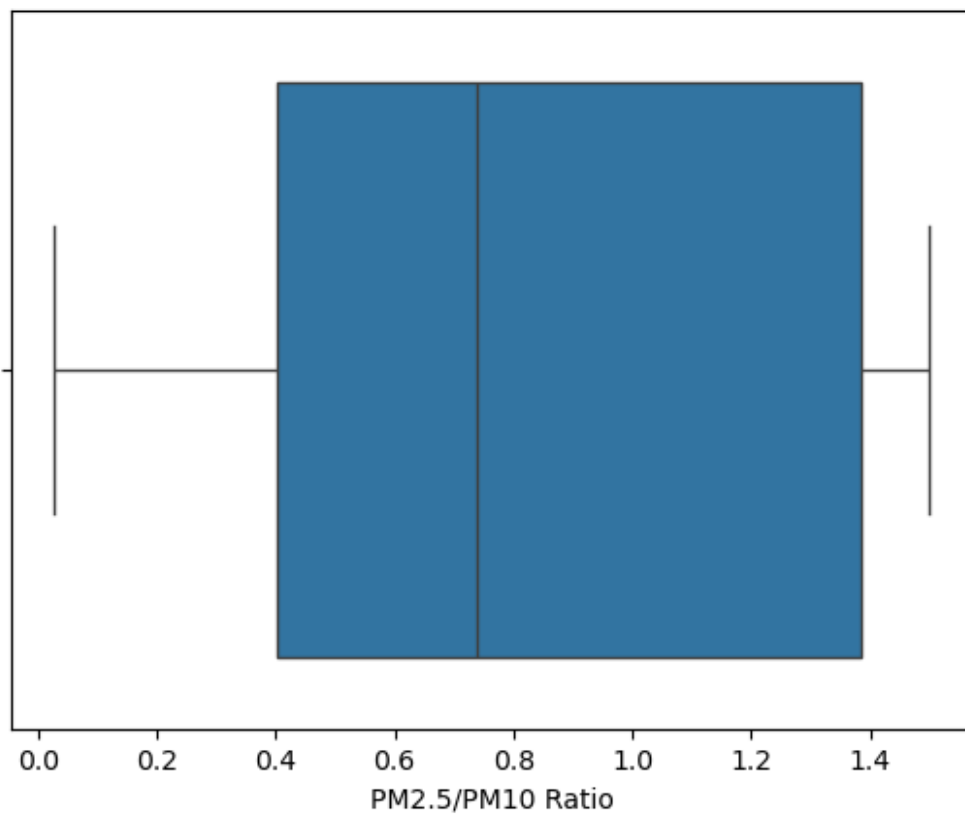
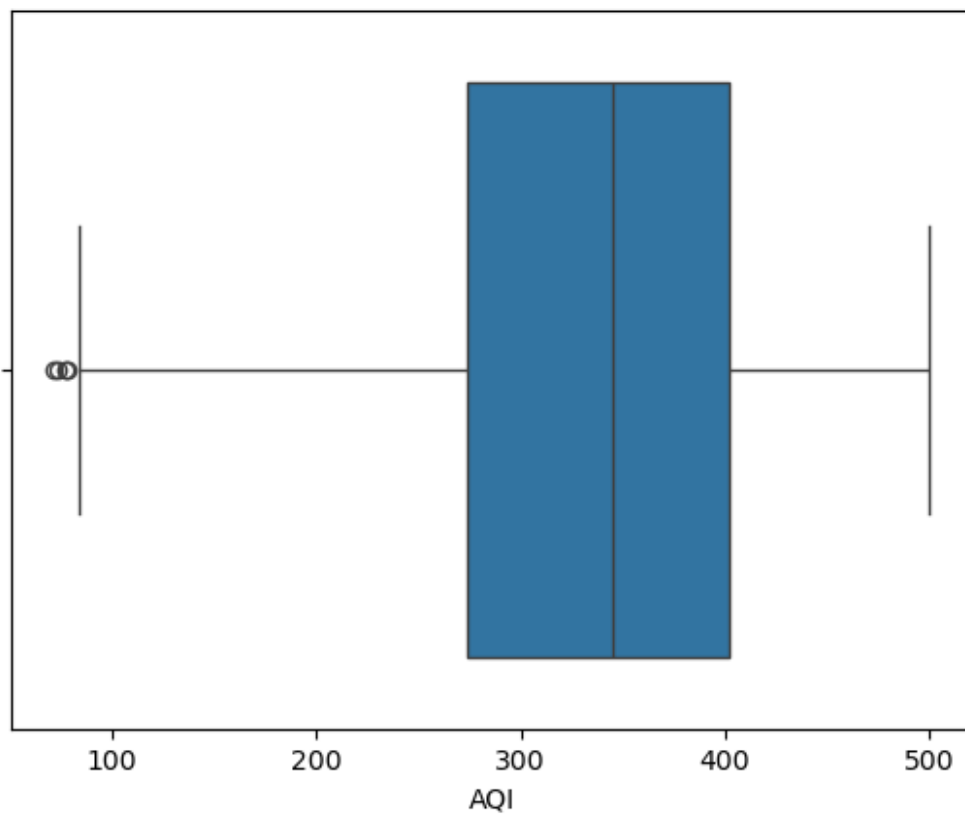
```
#Boxplot to indentify outliers  
#There is no outliers  
import warnings  
warnings.filterwarnings("ignore")  
for i in df.select_dtypes(include="number").columns:  
    sns.boxplot(data=df,x=i)  
    plt.show()
```





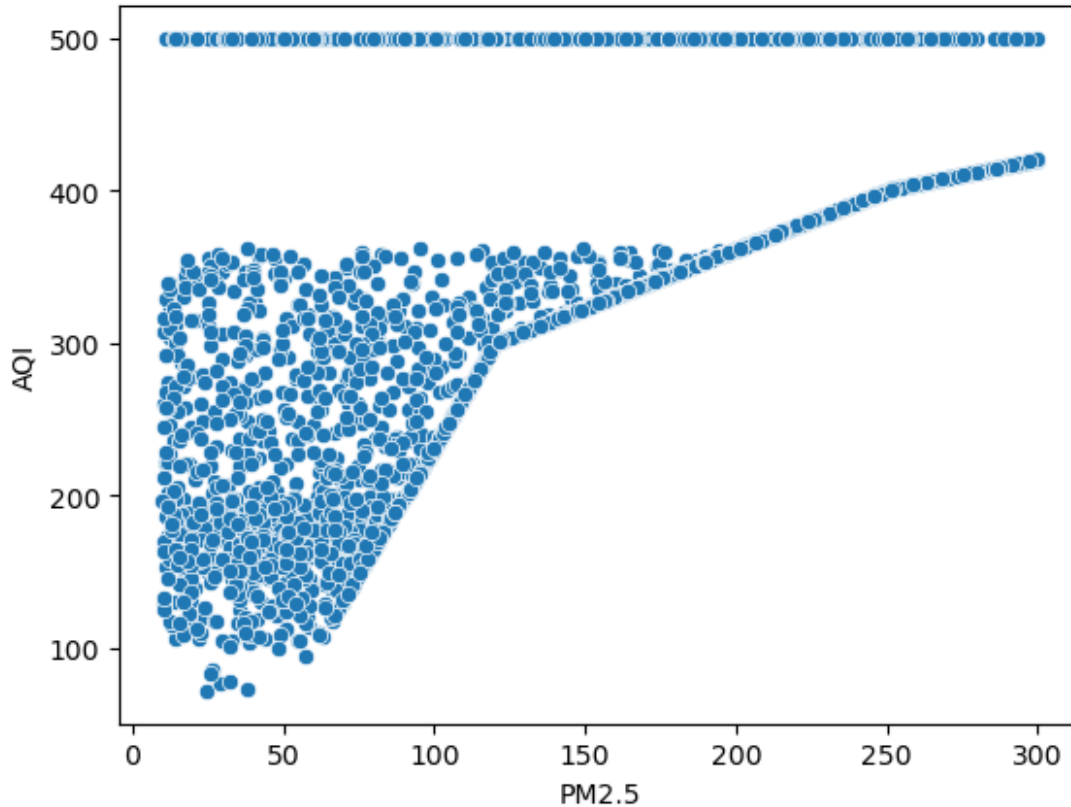


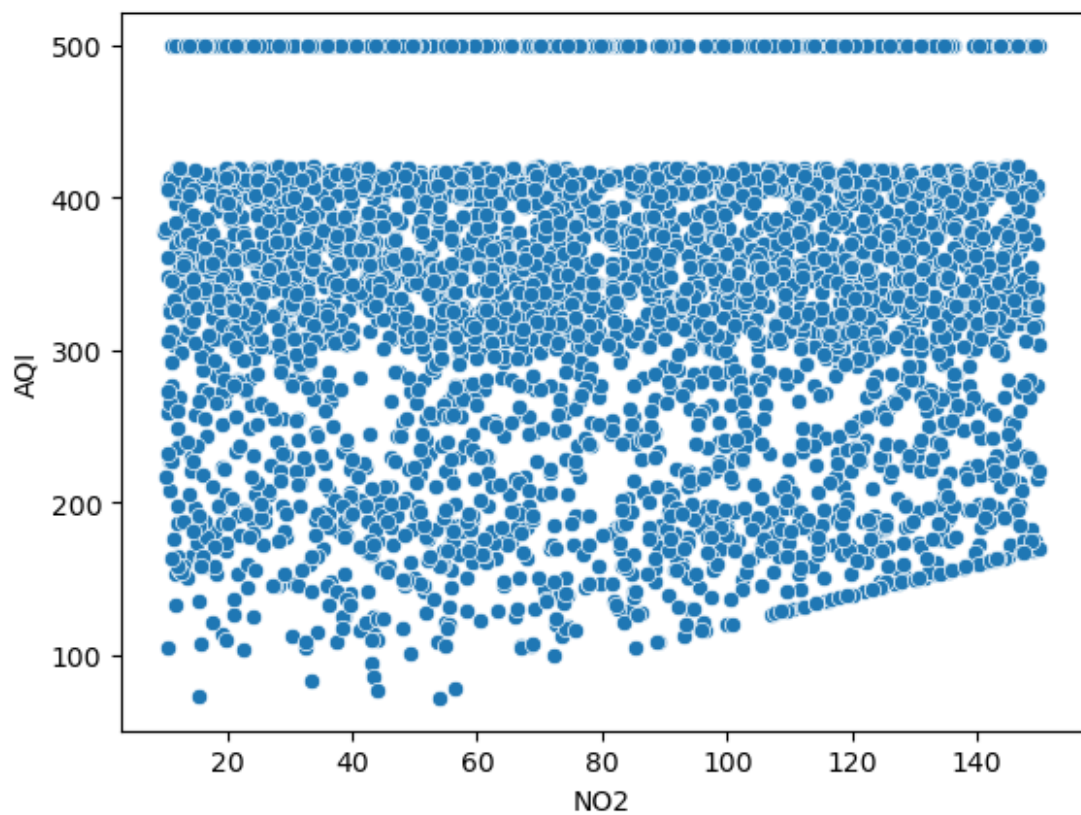
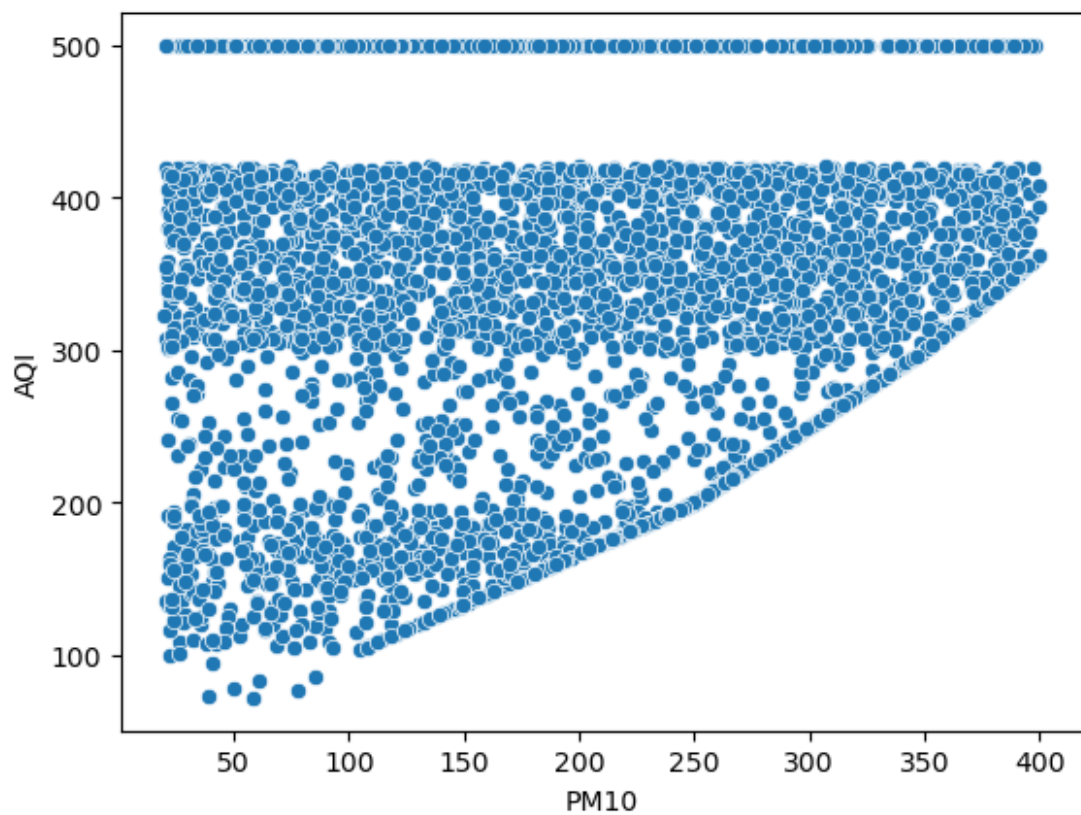


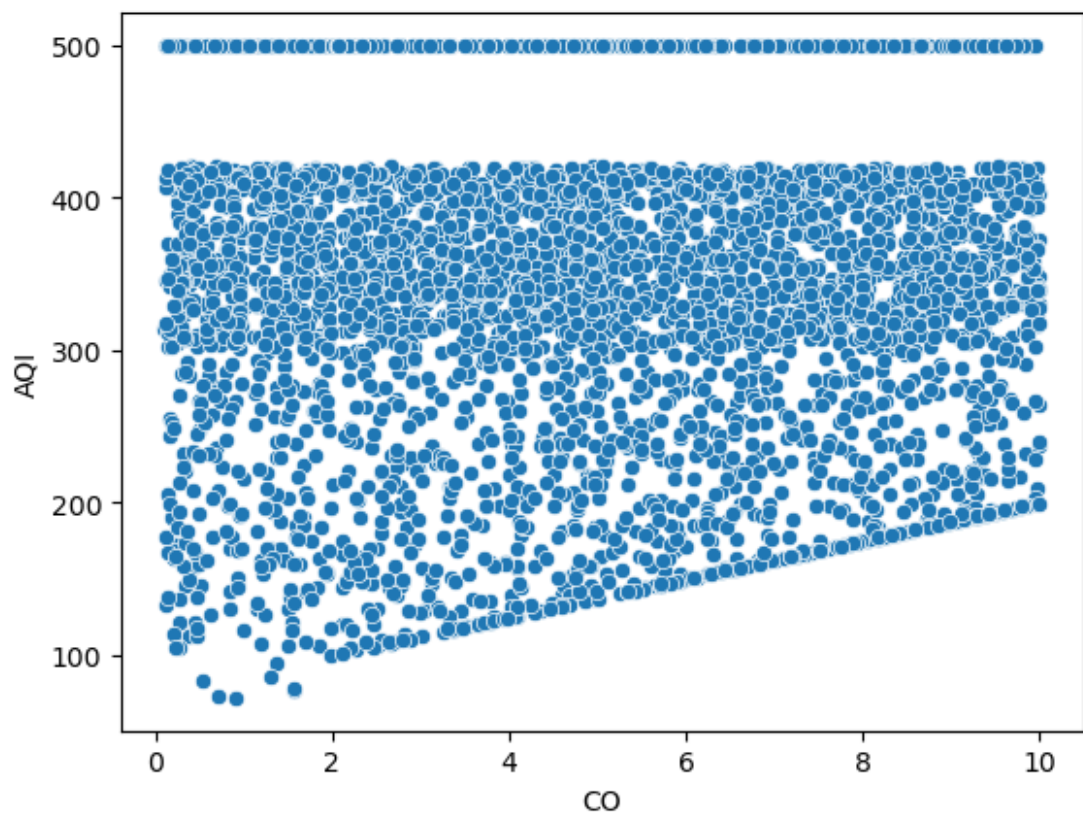
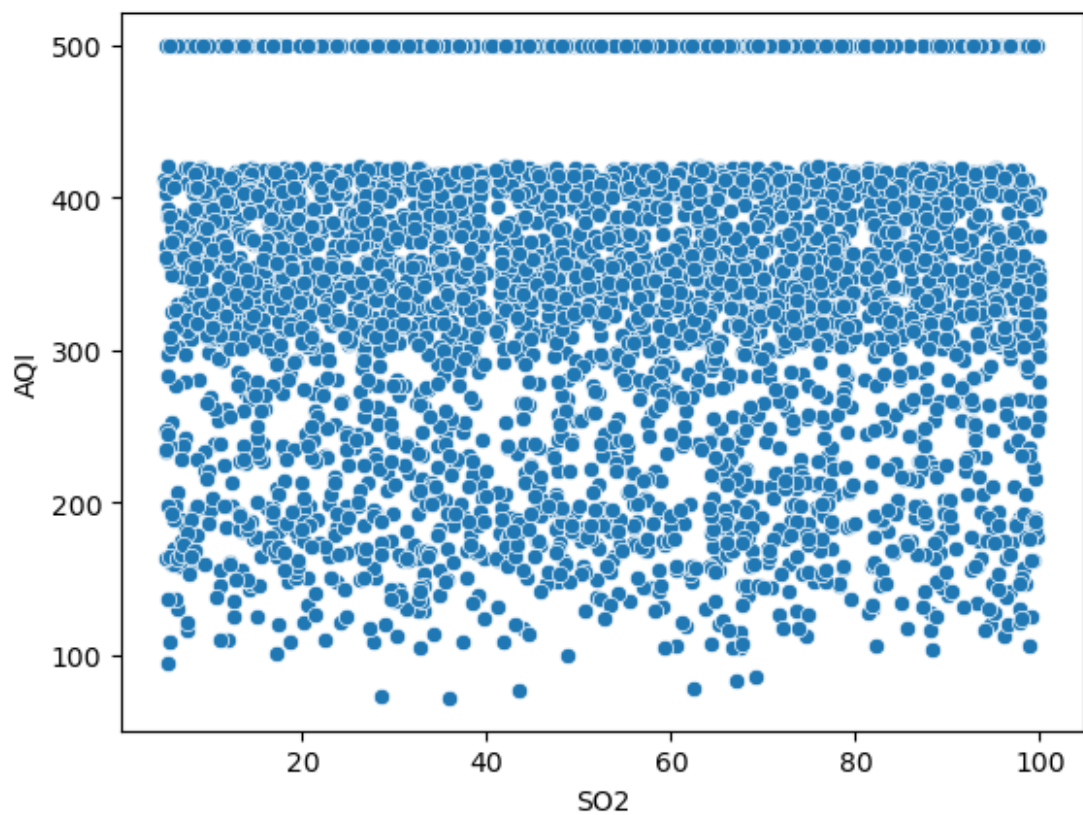


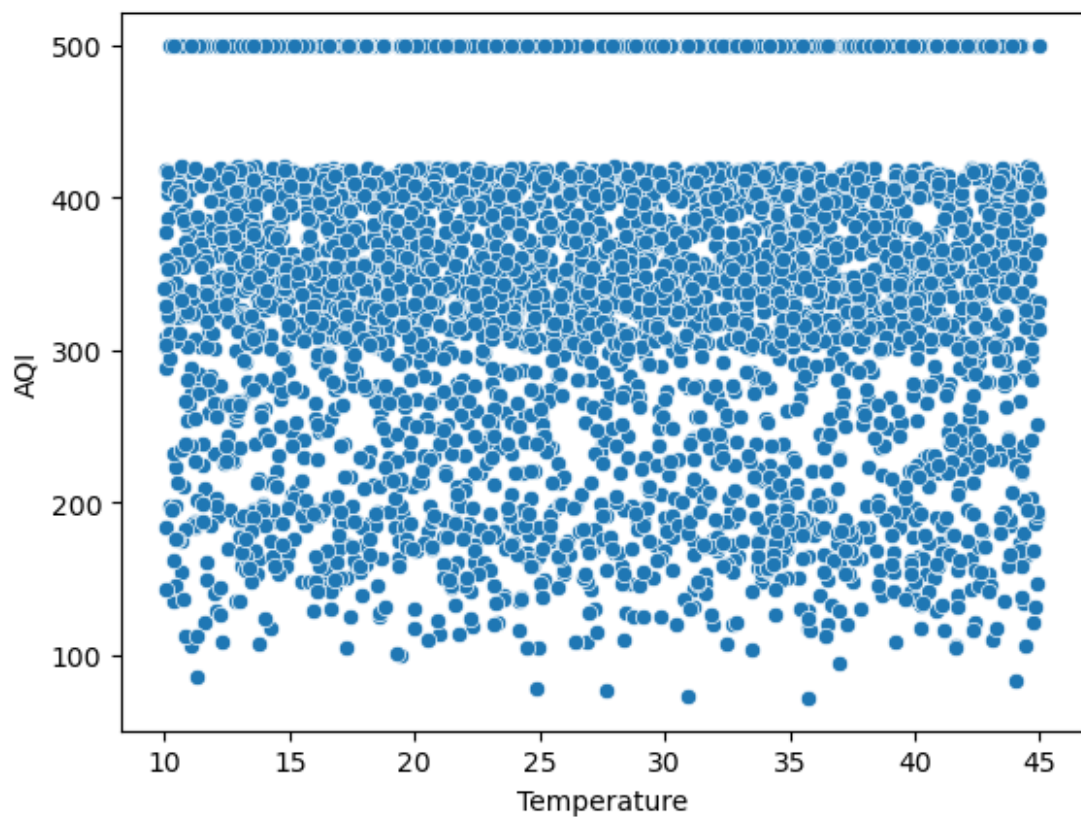
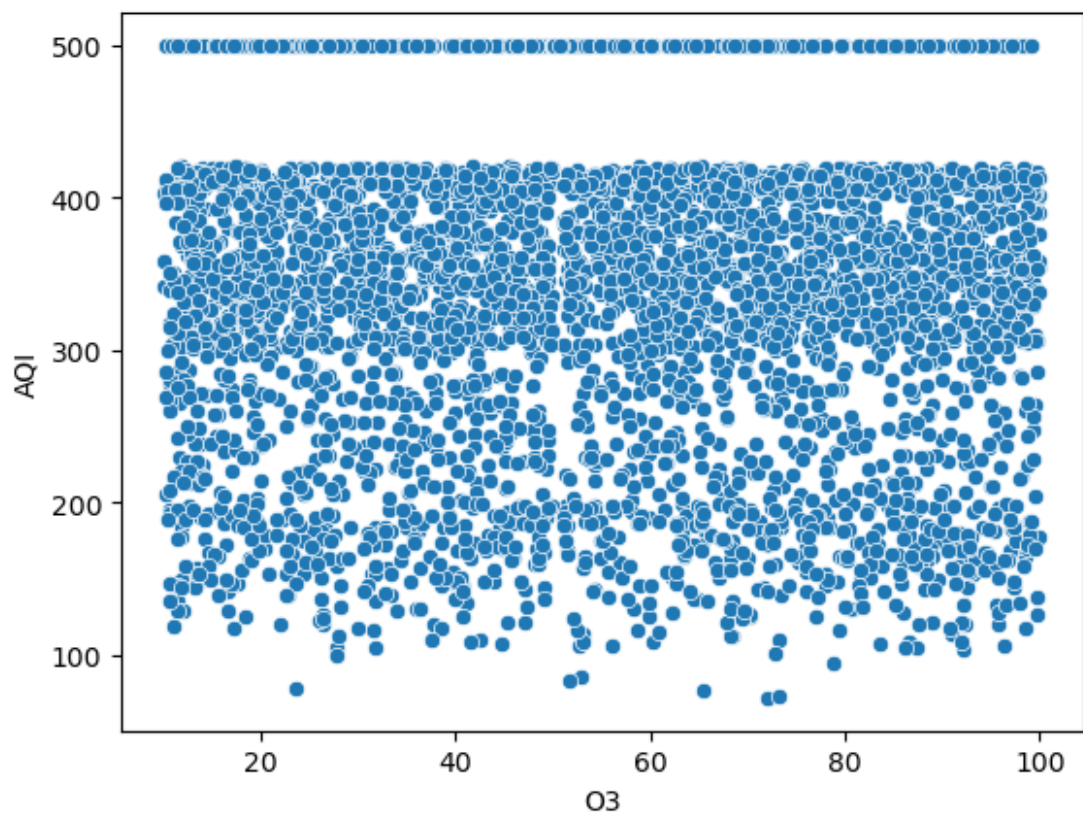

```
#Scatter Plot to understand the relationship
# df.select_dtypes(include="number").columns

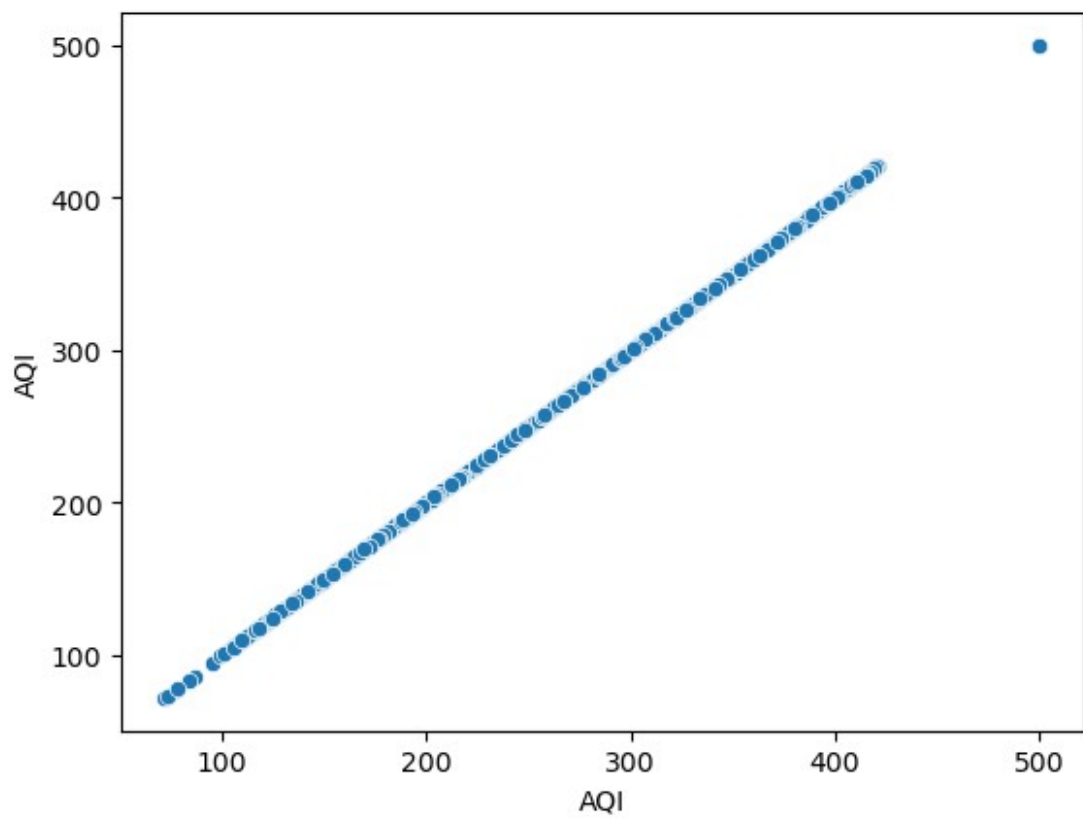
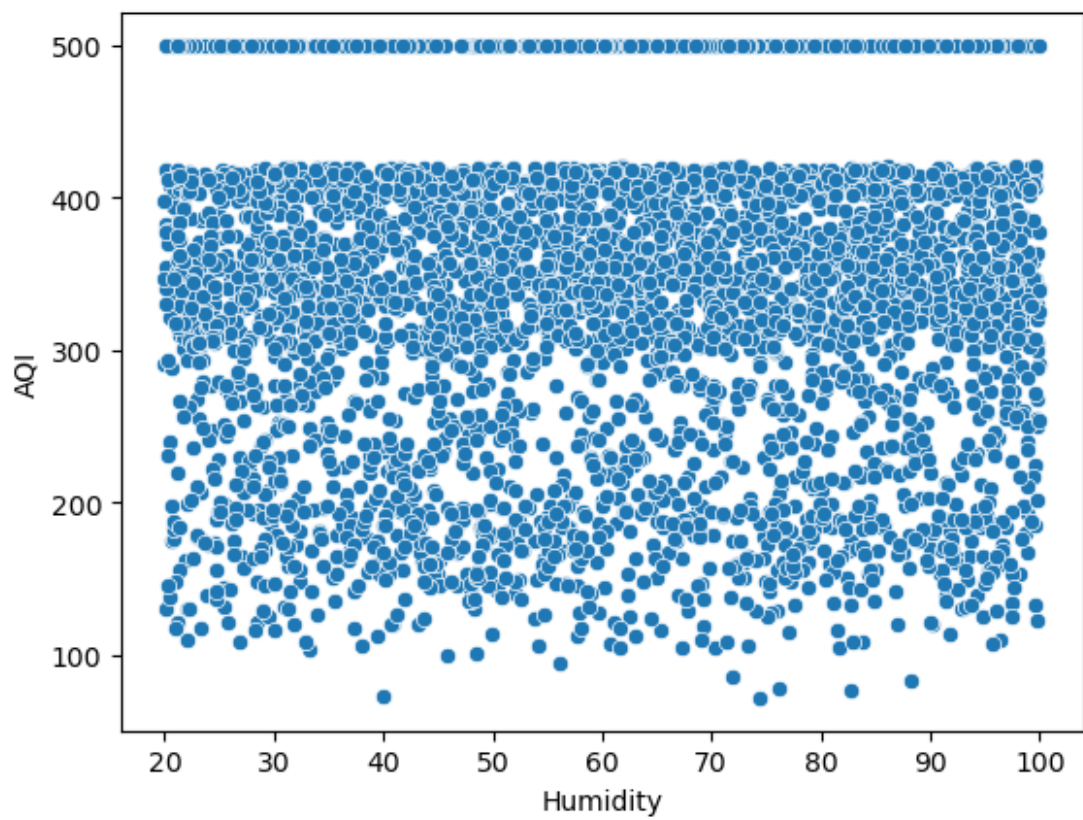
for i in ['PM2.5', 'PM10', 'N02', 'S02', 'CO', 'O3', 'Temperature',
'Humidity',
'AQI', 'PM2.5/PM10 Ratio']:
    sns.scatterplot(data=df, x=i, y='AQI')
plt.show()
```

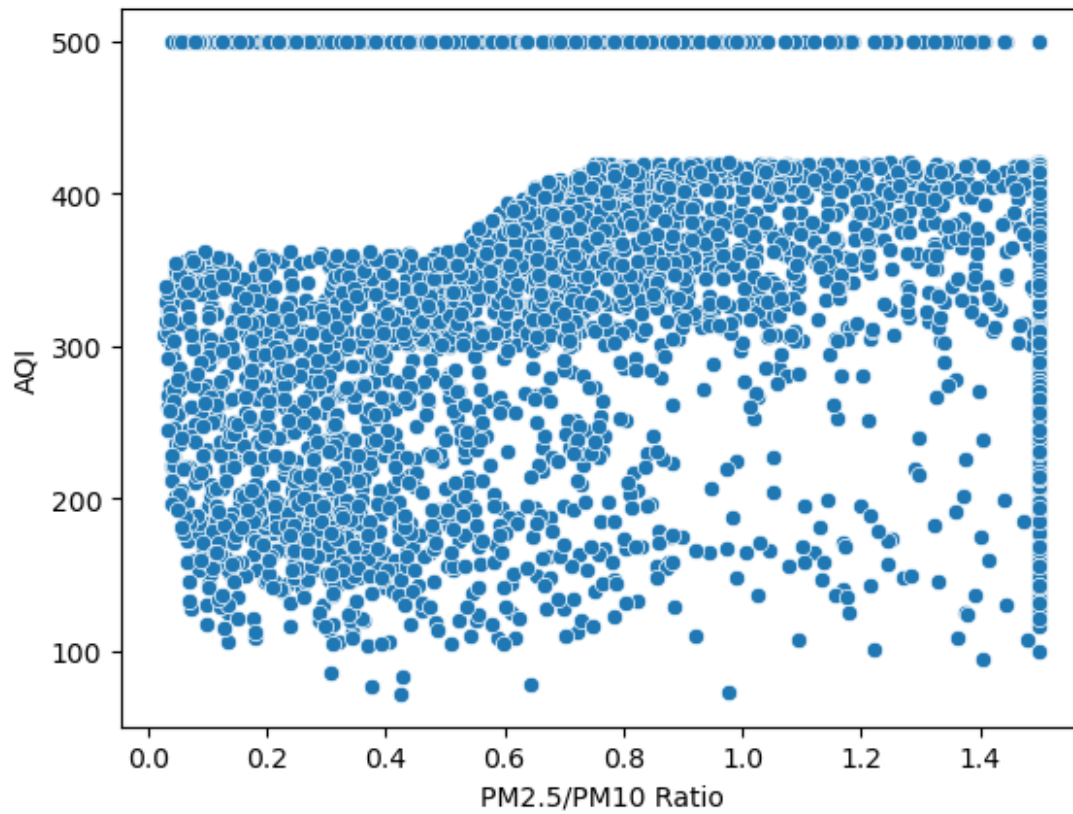












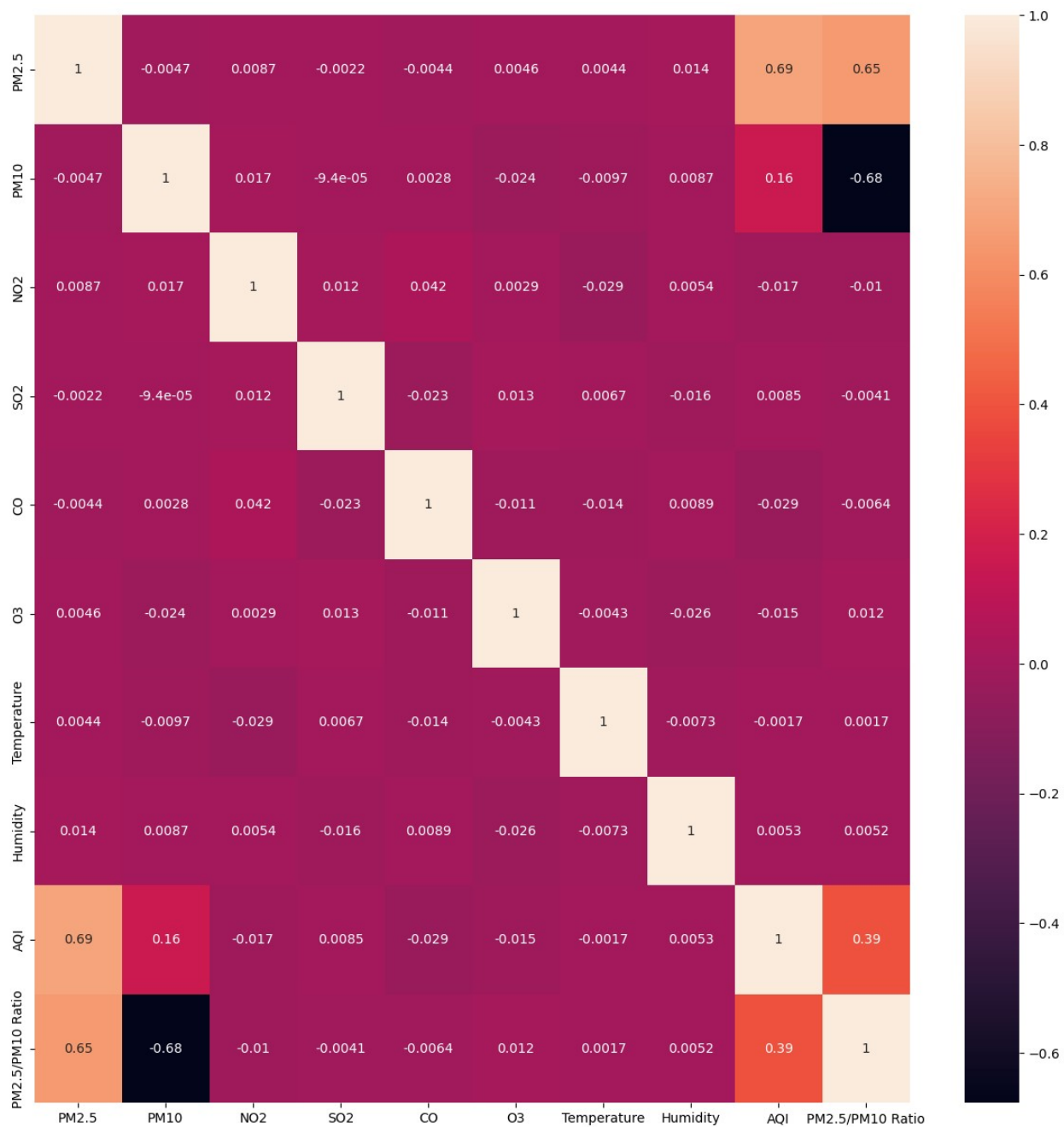
#correlation with heatmap to interpret the relation and multicolliniarity

```
s=df.select_dtypes(include="number").corr()
```

```
plt.figure(figsize=(15,15))
```

```
sns.heatmap(s,annot=True)
```

<Axes: >



STEP-5 MISSING VALUE TREATMENTS

Choose the method of imputing missing value

Like mean, median, mode or KNNImputer

There is no missing value in this dataset

STEP-6 OUTLIER TREATMENT

There is no outlier in this dataset

STEP-7 DUPLICATES AND GARBAGE VALUE TREATMENT

There is no duplicate value in this data but there are garbage values present in this dataset

To remove duplicates we will use:

`df.drop_duplicates()` function

To remove Garbage value , we will use the mean, median, mode function to change the value