

# Air Quality Index Prediction

## Problem Statement

Air pollution is being measured by AQI (Air Quality Index) and higher the AQI higher the air pollution. In this project we will analyze the air quality data of major developing cities in India to find some or patterns which will help us in predicting the AQI with the given data. Since this data has been captured at specific intervals over a period of time we will be conducting a various Analysis to predict the AQI of a city at a point in time.

## Importing Libraries that are useful for data manipulation and visualization

In [1]:

```
1  # For data reading / manipulation :
2  import pandas as pd
3
4  #For handling spatial data
5  import geopandas as gpd
6
7  #import psycopg2 library to connect database :
8  import psycopg2
9
10 # For reading the array :
11 import numpy as np
12
13 # For visualize the data and plotting :
14 import matplotlib.pyplot as plt
15 plt.rcParams['figure.figsize'] = [15,8]
16 import seaborn as sns
17
18 # To ignore warnings :
19 from warnings import filterwarnings
20 filterwarnings('ignore')
21
22 # Train-test split :
23 from sklearn.model_selection import train_test_split
24
25 #import LinearRegression from sklearn library
26 from sklearn.linear_model import LinearRegression
27
28 #import RandomForestRegressor from sklearn library
29 from sklearn.ensemble import RandomForestRegressor
30
31 from sklearn.tree import DecisionTreeRegressor
32
33 # To check the accuracy of model :
34 from sklearn.metrics import r2_score, mean_absolute_error
35
36 # Import 'stats' libraries for modeling :
37 from scipy import stats
38 import statsmodels.api as sm
39 from statsmodels.api import OLS
40
41 # To scaled the data :
42 from sklearn.preprocessing import StandardScaler
43
44
```

## Load and Read the Data from Database

In [3]:

```

1 #create a connection object
2 connection=psycopg2.connect(database="library",user="postgres",password="snag2568")
3
4 #create a cursor
5 cursor=connection.cursor()
6 connection.commit()
7 AQI_city_day=pd.read_sql("SELECT * FROM raw_data",con=connection)
8

```

In [4]:

```

1 # Read the data :
2 AQI_city_day=pd.read_sql("SELECT * FROM raw_data",con=connection)
3
4
5
6 # Print the first five observations to see the data structure :
7 AQI_city_day.head()

```

Out[4]:

	City	Date	PM2_5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3	Benzene
0	Ahmedabad	2015-01-01	NaN	NaN	0.92	18.22	17.15	NaN	0.92	27.64	133.36	0.01
1	Ahmedabad	2015-01-02	NaN	NaN	0.97	15.69	16.46	NaN	0.97	24.55	34.06	3.61
2	Ahmedabad	2015-01-03	NaN	NaN	17.40	19.30	29.70	NaN	17.40	29.07	30.70	6.81
3	Ahmedabad	2015-01-04	NaN	NaN	1.70	18.48	17.97	NaN	1.70	18.59	36.08	4.41
4	Ahmedabad	2015-01-05	NaN	NaN	22.10	21.42	37.76	NaN	22.10	39.33	39.31	7.01

In [5]:

```

1 geo=gpd.GeoDataFrame(AQI_city_day,geometry=gpd.points_from_xy(AQI_city_day.Latitude,

```

In [6]:

```
1 geo.head(3)
```

Out[6]:

	City	Date	PM2_5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3	Benzene
0	Ahmedabad	2015-01-01	NaN	NaN	0.92	18.22	17.15	NaN	0.92	27.64	133.36	0.01
1	Ahmedabad	2015-01-02	NaN	NaN	0.97	15.69	16.46	NaN	0.97	24.55	34.06	3.61
2	Ahmedabad	2015-01-03	NaN	NaN	17.40	19.30	29.70	NaN	17.40	29.07	30.70	6.81

## Understand the Data

In [8]:

```
1 # Checking shape and data types of the data :
2
3 AQI_city_day.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 29531 entries, 0 to 29530
Data columns (total 18 columns):
#   Column      Non-Null Count  Dtype
---  -
0   City         29531 non-null  object
1   Date         29531 non-null  object
2   PM2_5        24933 non-null  float64
3   PM10         18391 non-null  float64
4   NO           25949 non-null  float64
5   NO2          25946 non-null  float64
6   NOx          25346 non-null  float64
7   NH3          19203 non-null  float64
8   CO           27472 non-null  float64
9   SO2          25677 non-null  float64
10  O3           25509 non-null  float64
11  Benzene      23908 non-null  float64
12  Toluene      21490 non-null  float64
13  Xylene       11422 non-null  float64
14  AQI          24850 non-null  float64
15  Longitude    29531 non-null  float64
16  Latitude     29531 non-null  float64
17  geometry     29531 non-null  geometry
dtypes: float64(15), geometry(1), object(2)
memory usage: 4.1+ MB
```

In [9]:

```
1 # Summary Statistics of the numeric data :  
2  
3 AQI_city_day.describe()
```

Out[9]:

	PM2_5	PM10	NO	NO2	NOx	NH3
count	24933.000000	18391.000000	25949.000000	25946.000000	25346.000000	19203.000000
mean	67.450578	118.127103	17.574730	28.560659	32.309123	23.483476
std	64.661449	90.605110	22.785846	24.474746	31.646011	25.684275
min	0.040000	0.010000	0.020000	0.010000	0.000000	0.010000
25%	28.820000	56.255000	5.630000	11.750000	12.820000	8.580000
50%	48.570000	95.680000	9.890000	21.690000	23.520000	15.850000
75%	80.590000	149.745000	19.950000	37.620000	40.127500	30.020000
max	949.990000	1000.000000	390.680000	362.210000	467.630000	352.890000

### Interpretation

- Summary statistics give maximum, minimum, average, first and third quartile and standard deviation of the data.
- Most of the variables having minimum value as 0.

## Data preparation and cleaning

In [10]:

```
1 # Checking for the missing values and its precentage :
2
3 values = AQI_city_day.isnull().sum()
4 percentage = 100*AQI_city_day.isnull().sum()/len(AQI_city_day)
5 table = pd.concat([values,percentage.round(2)],axis=1)
6 table.columns = ['No of missing values','Percent of missing values']
7 table[table['No of missing values']!=0].sort_values('Percent of missing values',asce
```

Out[10]:

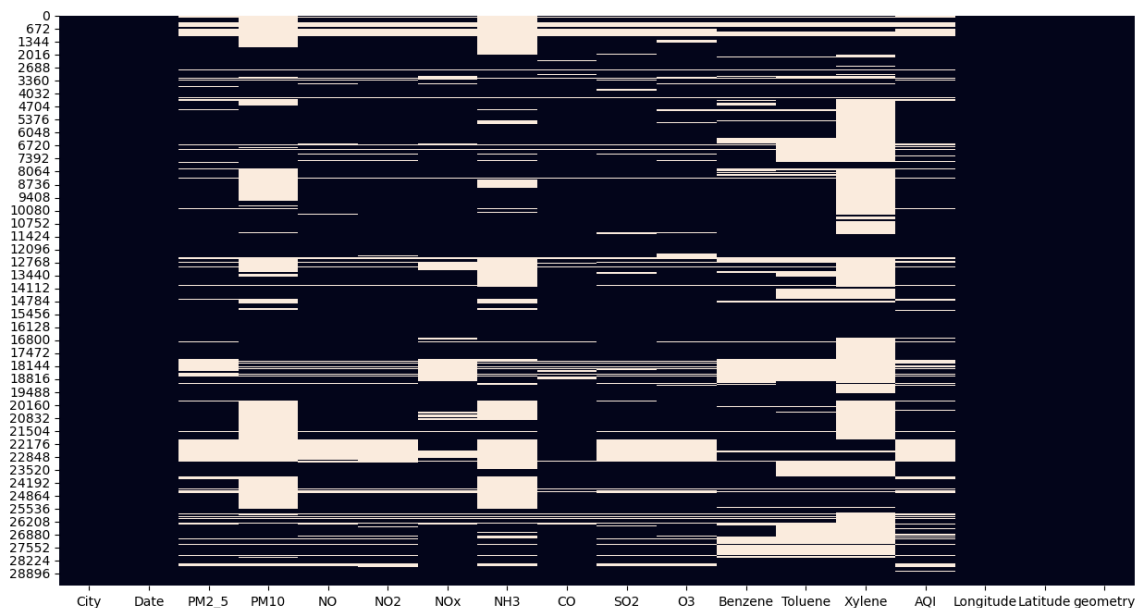
	No of missing values	Percent of missing values
Xylene	18109	61.320000
PM10	11140	37.720000
NH3	10328	34.970000
Toluene	8041	27.230000
Benzene	5623	19.040000
AQI	4681	15.850000
PM2_5	4598	15.570000
NOx	4185	14.170000
O3	4022	13.620000
SO2	3854	13.050000
NO2	3585	12.140000
NO	3582	12.130000
CO	2059	6.970000

In [11]:

```

1 # Check the missing values with heatmap :
2
3 sns.heatmap(AQI_city_day.isnull(), cbar=False)
4 plt.show()
5

```



In [12]:

```

1 # Try methods to impute missing values :
2
3 AQI_city_day.groupby(['City', 'Date'])['AQI'].mean()
4

```

Out[12]:

City	Date	AQI
Ahmedabad	2015-01-01	NaN
	2015-01-02	NaN
	2015-01-03	NaN
	2015-01-04	NaN
	2015-01-05	NaN
		...
Visakhapatnam	2020-06-27	41.0
	2020-06-28	70.0
	2020-06-29	68.0
	2020-06-30	54.0
	2020-07-01	50.0

Name: AQI, Length: 29531, dtype: float64

In [13]:

```

1 # Converting Datetime variable into datetime data type :
2
3 AQI_city_day['Date'] = AQI_city_day['Date'].apply(pd.to_datetime)

```

In [14]:

```
1 # Check whether it is converted or not :  
2  
3 AQI_city_day.dtypes
```

Out[14]:

```
City                object  
Date               datetime64[ns]  
PM2_5              float64  
PM10               float64  
NO                 float64  
NO2                float64  
NOx                float64  
NH3                float64  
CO                 float64  
SO2                float64  
O3                 float64  
Benzene            float64  
Toluene            float64  
Xylene             float64  
AQI                float64  
Longitude           float64  
Latitude            float64  
geometry            geometry  
dtype: object
```

- Datetime variable is changed into datetime datatype.

In [15]:

```
1 # Impute the missing values by grouping city with and average of 5 days:  
2  
3 data_fill = AQI_city_day.iloc[:, :15].fillna(AQI_city_day.iloc[:, :15].groupby(['Cit
```



In [16]:

```

1 # After inputing missing values by day checking for the missing values :
2
3 values = data_fill.isnull().sum()
4 percentage = 100*data_fill.isnull().sum()/len(data_fill)
5 table = pd.concat([values,percentage.round(2)],axis=1)
6 table.columns = ['No of missing values','Percent of missing values']
7 table[table['No of missing values']!=0].sort_values('Percent of missing values',asce

```

Out[16]:

	No of missing values	Percent of missing values
<b>Xylene</b>	17618	59.660000
<b>PM10</b>	10495	35.540000
<b>NH3</b>	9742	32.990000
<b>Toluene</b>	7396	25.040000
<b>Benzene</b>	4888	16.550000
<b>PM2_5</b>	3990	13.510000
<b>NOx</b>	3700	12.530000
<b>AQI</b>	3661	12.400000
<b>O3</b>	3201	10.840000
<b>SO2</b>	3050	10.330000
<b>NO2</b>	2885	9.770000
<b>NO</b>	2878	9.750000
<b>CO</b>	1354	4.590000

- Still the missing values are present in the data, let's try to fill with average of month.

In [17]:

```

1 # Impute the missing values by grouping city and month :
2
3 data_fill = data_fill.fillna(data_fill.groupby(['City', pd.Grouper(key='Date', freq=

```

In [18]:

```

1 # After inputing missing values by month check for the missing values :
2
3 values = data_fill.isnull().sum()
4 percentage = 100*data_fill.isnull().sum()/len(data_fill)
5 table = pd.concat([values,percentage.round(2)],axis=1)
6 table.columns = ['No of missing values','Percent of missing values']
7 table[table['No of missing values']!=0].sort_values('Percent of missing values',asce

```

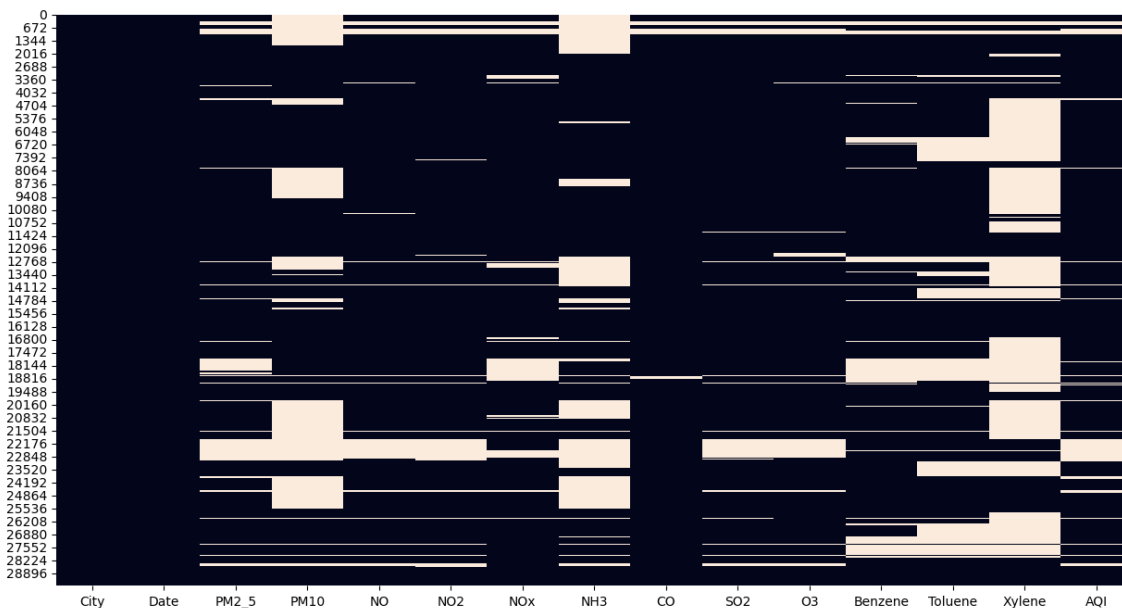
	No of missing values	Percent of missing values
Xylene	17033	57.680000
PM10	9570	32.410000
NH3	8900	30.140000
Toluene	6742	22.830000
Benzene	3999	13.540000
PM2_5	3067	10.390000
NOx	3001	10.160000
AQI	2596	8.790000
O3	2113	7.160000
NO2	2056	6.960000
NO	1964	6.650000
SO2	1935	6.550000

In [19]:

```

1 sns.heatmap(data_fill.isnull(), cbar=False)
2 plt.show()

```



In [20]:

```

1 # Imputing missing values by backward fill :
2
3 data_fill = data_fill.fillna(method = 'bfill',axis=0)

```

In [21]:

```

1 # Still data is missing in the columns, hence using forward fill to impute this :
2
3 data_fill = data_fill.fillna(method = 'ffill',axis=0)

```

In [22]:

```

1 # After inputing missing values by bfill abd ffill:
2
3 values = data_fill.isnull().sum()
4 percentage = 100*data_fill.isnull().sum()/len(data_fill)
5 table = pd.concat([values,percentage.round(2)],axis=1)
6 table.columns = ['No of missing values','Percent of missing values']
7 table[table['No of missing values']!=0].sort_values('Percent of missing values',asce

```

Out[22]:

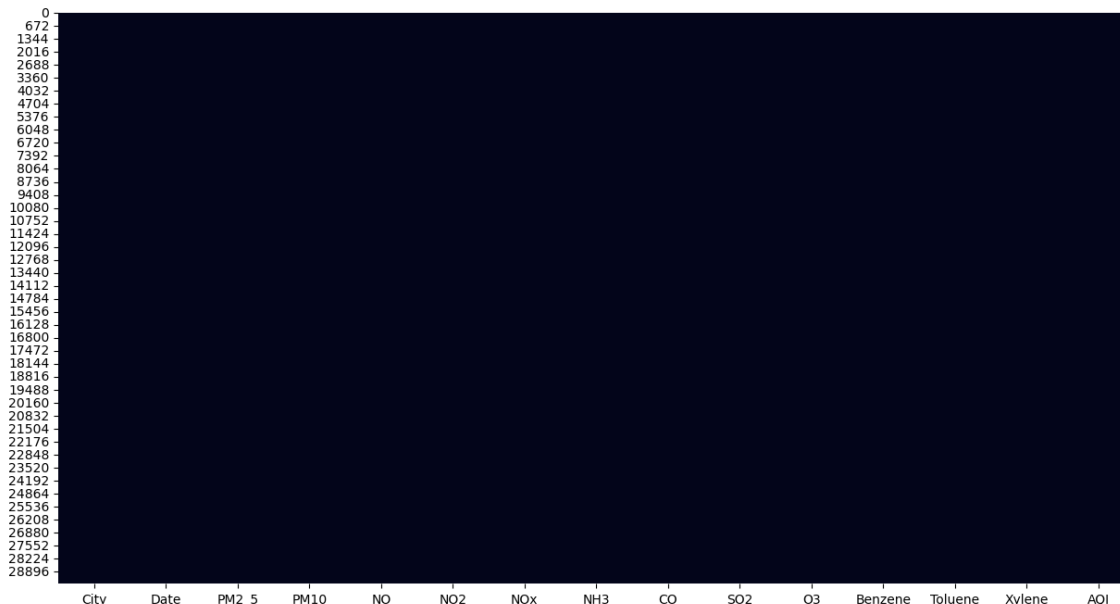
No of missing values	Percent of missing values
----------------------	---------------------------

In [23]:

```

1 sns.heatmap(data_fill.isnull(), cbar=False)
2 plt.show()

```



In [24]:

```
1 print(data_fill.columns)
```

```
Index(['City', 'Date', 'PM2_5', 'PM10', 'NO', 'NO2', 'NOx', 'NH3', 'CO',  
'SO2',  
      'O3', 'Benzene', 'Toluene', 'Xylene', 'AQI'],  
      dtype='object')
```

In [25]:

```
1 data_fill.head()
```

Out[25]:

	City	Date	PM2_5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3
0	Ahmedabad	2015-01-01	81.367222	137.342727	0.92	18.22	17.15	26.64	0.92	27.64	133.36
1	Ahmedabad	2015-01-02	81.367222	137.342727	0.97	15.69	16.46	26.64	0.97	24.55	34.06
2	Ahmedabad	2015-01-03	81.367222	137.342727	17.40	19.30	29.70	26.64	17.40	29.07	30.70
3	Ahmedabad	2015-01-04	81.367222	137.342727	1.70	18.48	17.97	26.64	1.70	18.59	36.08
4	Ahmedabad	2015-01-05	81.367222	137.342727	22.10	21.42	37.76	26.64	22.10	39.33	39.31

In [26]:

```
1 # After imputing missing values, check summary statistics of the data :  
2  
3 data_fill.describe()
```

Out[26]:

	PM2_5	PM10	NO	NO2	NOx	NH3
count	29531.000000	29531.000000	29531.000000	29531.000000	29531.000000	29531.000000
mean	67.612826	101.114376	21.370930	27.232570	35.339042	21.864417
std	64.021390	84.360941	28.809299	24.026101	33.420449	23.376915
min	0.040000	0.010000	0.020000	0.010000	0.000000	0.010000
25%	26.992500	44.560000	5.900000	10.610000	13.885000	6.756200
50%	47.900000	72.040000	10.620000	20.620000	24.590000	15.340000
75%	84.650000	137.342727	23.255000	36.420000	44.087111	29.430000
max	949.990000	1000.000000	390.680000	362.210000	467.630000	352.890000

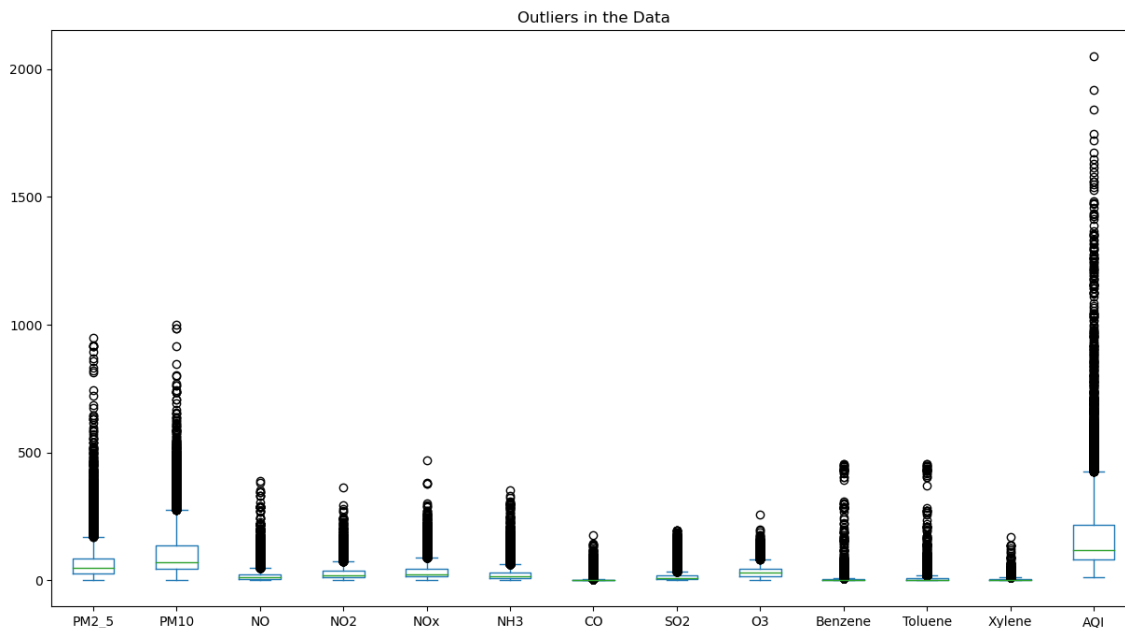
- After imputing missing values by this method, mean and the standard deviation is not much changed.

In [27]:

```
1 data=data_fill.copy()
```

In [28]:

```
1 # Checking for outliers :
2
3 data_fill.plot(kind='box')
4 plt.title("Outliers in the Data")
5 plt.show()
6 plt.savefig('Outliers.png', dpi=300, bbox_inches='tight')
```



&lt;Figure size 1500x800 with 0 Axes&gt;

In [29]:

```
1 from scipy import stats
2
3 # Assuming your DataFrame is called "df"
4 # Calculate the z-scores for each column
5 z_scores = stats.zscore(data[['PM2_5', 'PM10', 'NO', 'NO2', 'NOx', 'NH3', 'CO', 'SO2', 'O3', 'Benzene', 'Toluene', 'Xylene', 'AQI']])
6
7 # Set a threshold z-score above which a data point is considered an outlier
8 threshold = 3
9
10 # Remove outliers by filtering out all rows with a z-score greater than the threshold
11 data = data[(z_scores < threshold).all(axis=1)]
12
```

In [30]:

```
1 data.isnull().sum()
```

Out[30]:

```
City          0
Date          0
PM2_5        614
PM10         685
NO           428
NO2          524
NOx          467
NH3          459
CO           548
SO2         1441
O3           399
Benzene       158
Toluene       319
Xylene        263
AQI          340
dtype: int64
```

In [31]:

```
1 data.dropna(inplace=True)
```

In [32]:

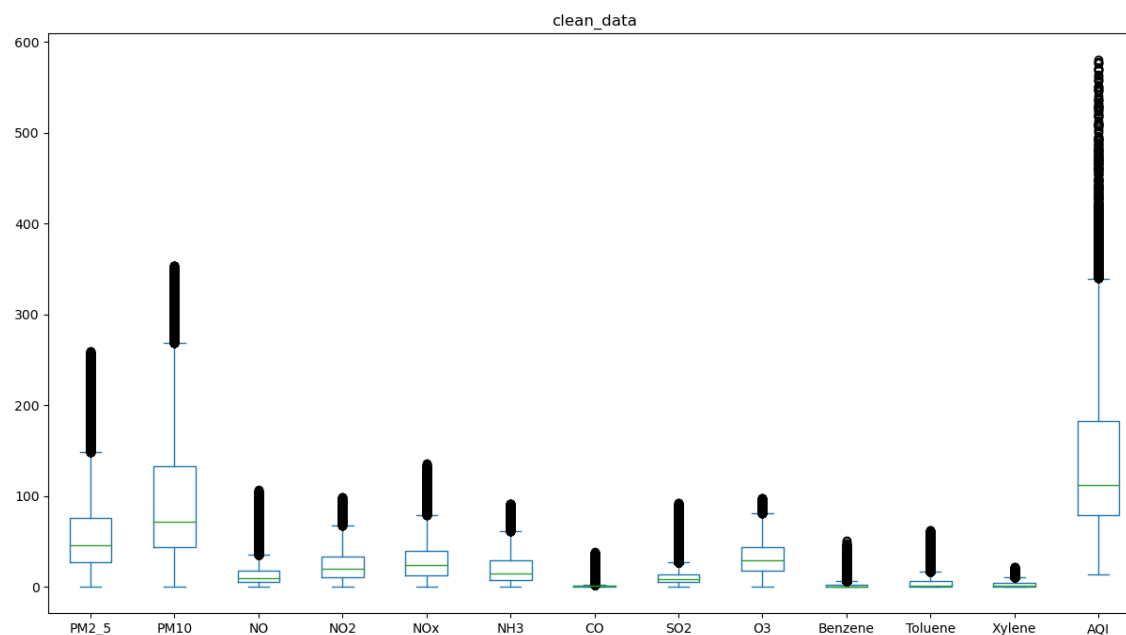
```
1 data.isnull().sum()
```

Out[32]:

```
City          0
Date          0
PM2_5          0
PM10           0
NO             0
NO2            0
NOx            0
NH3            0
CO             0
SO2            0
O3             0
Benzene        0
Toluene        0
Xylene         0
AQI            0
dtype: int64
```

In [33]:

```
1 data.plot(kind='box')
2 plt.title("clean_data")
3 plt.show()
4 plt.savefig('Outliers.png', dpi=300, bbox_inches='tight')
```



&lt;Figure size 1500x800 with 0 Axes&gt;

## Interpretation

- From the box plot, can say that the outliers are present in the data, but this is the time series data.
- So we can not remove all outliers, if we do this the continuity will not remain for the data.

# Univariate ad Multivariate Month wise Analysis

In [34]:

```

1  # In this we are not considering the AQI_Bucket column, because it is not nessasary
2  # Here, we do analysis on all over country, so Aggregating the data into month wise
3
4  AQI_df = data.groupby(['City', (data.Date.dt.strftime('%Y-%m'))]).mean()
5
6  # Reset index :
7  AQI_df = AQI_df.reset_index()
8
9  # Check the First 5 observation :
10 AQI_df.head()

```

Out[34]:

	City	Date	PM2_5	PM10	NO	NO2	NOx	NH3	CC
0	Ahmedabad	2015-01	80.740899	137.342727	8.891429	23.060000	28.761429	26.64	8.891429
1	Ahmedabad	2015-02	83.155000	137.342727	8.047105	24.147368	30.995526	26.64	8.047105
2	Ahmedabad	2015-03	92.188333	137.342727	8.362262	24.424881	32.181905	26.64	8.362262
3	Ahmedabad	2015-04	98.601923	137.342727	6.158846	19.450769	24.903462	26.64	6.158846
4	Ahmedabad	2015-05	71.246552	137.342727	7.844828	16.702759	23.523448	26.64	7.844828

In [35]:

```
1 AQI_df.describe()
```

Out[35]:

	PM2_5	PM10	NO	NO2	NOx	NH3	CO
count	949.000000	949.000000	949.000000	949.000000	949.000000	949.000000	949.000000
mean	61.954768	95.400427	15.508434	25.523973	31.852424	20.595817	1.767571
std	42.660749	63.206121	13.600133	16.940917	25.032708	15.863441	3.346783
min	2.000000	3.555287	0.022500	0.070000	0.000000	0.017500	0.000000
25%	31.766429	47.514000	6.798333	11.837667	14.640948	7.883548	0.576667
50%	49.706957	79.441552	10.981429	22.331613	25.267500	16.240000	0.943704
75%	80.816667	137.342727	20.009000	34.977931	40.547619	30.046452	1.487200
max	221.335714	327.186000	90.910000	89.765000	121.174000	86.210000	30.282800

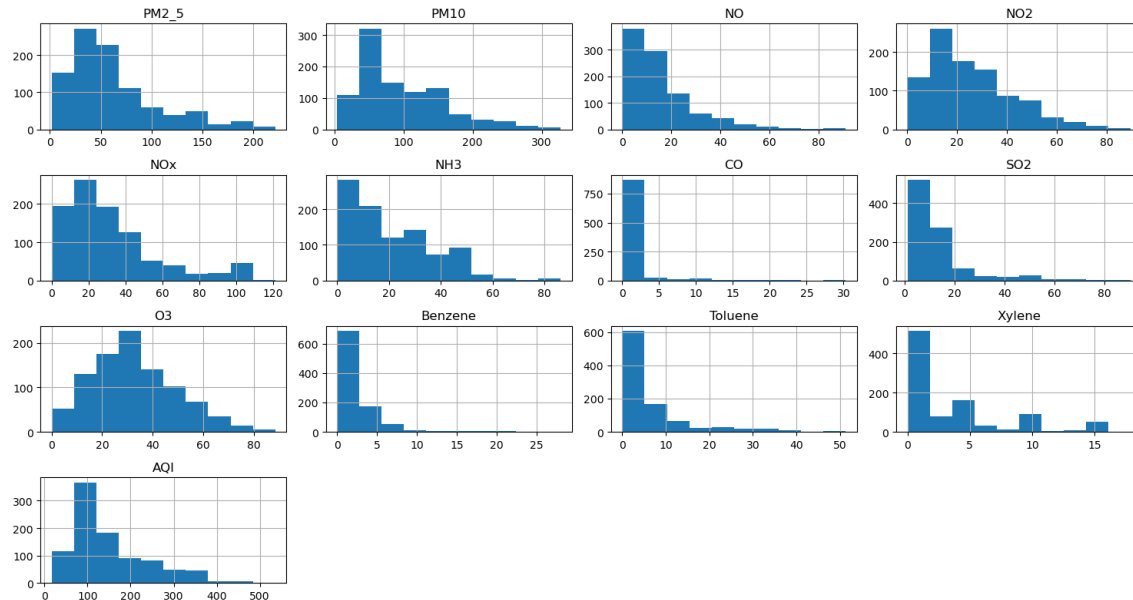


In [36]:

```

1 # Check the distribution of the all numerical columns and print skewness of the data
2
3 AQI_df.drop(['City', 'Date'], axis=1).hist()
4 plt.tight_layout()
5 plt.title("Distribution of Data")
6 plt.show()
7 plt.savefig('Distribution of data.png')
8 # Print the skewness of the data :
9 print('Skewness :\n', AQI_df.drop(['City', 'Date'], axis=1).skew())

```



Skewness :

```

PM2_5      1.310032
PM10       1.010859
NO         1.976403
NO2        0.912806
NOx        1.320127
NH3        0.892725
CO         4.841186
SO2        2.626591
O3         0.511202
Benzene    3.463085
Toluene    2.159469
Xylene     1.465082
AQI        1.212010
dtype: float64

```

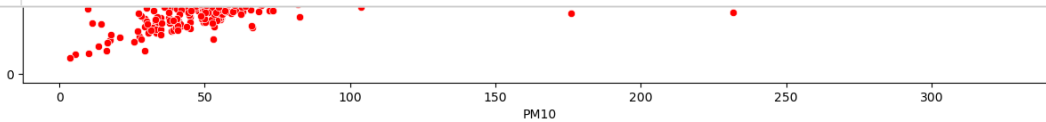
&lt;Figure size 1500x800 with 0 Axes&gt;

## Interpretation:

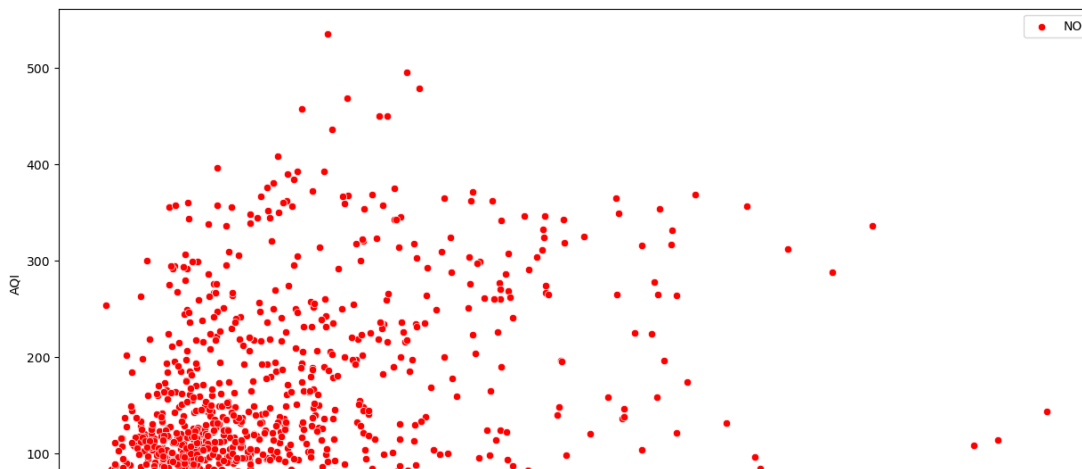
- From the plot and skewness, can say that the Benzene, Toluene, CO, NO2 are highly skewed and right skewed.

In [37]:

```
1 # Check the effect of all pollutants on AQI :
2
3 for i in AQI_df.iloc[:, 2:13]:
4     print('The impact of ', i, 'on AQI')
5     sns.scatterplot(x = i, y = 'AQI', data = AQI_df, marker="o", sizes=200, color="r", la
6     plt.legend()
7     plt.show()
8
```



The impact of NO on AQI



## Interpretations:

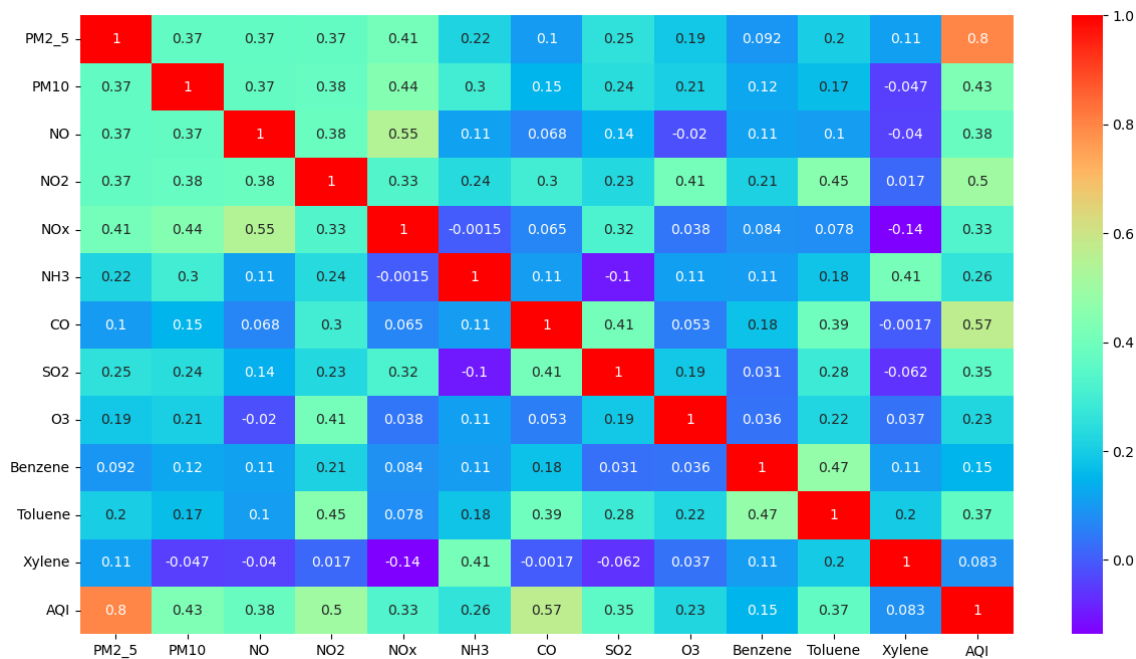
- From the scatterplot, can say that the PM2.5, PM10, CO has positive correlation with AQI.
- It means, when pollutants increase, the AQI increases.
- We will find the relationship of other variables using heatmap

In [38]:

```

1 # Correlation of the numerical data with heatmap :
2
3 sns.heatmap(AQI_df.corr(), annot=True, cmap='rainbow')
4 plt.show()

```

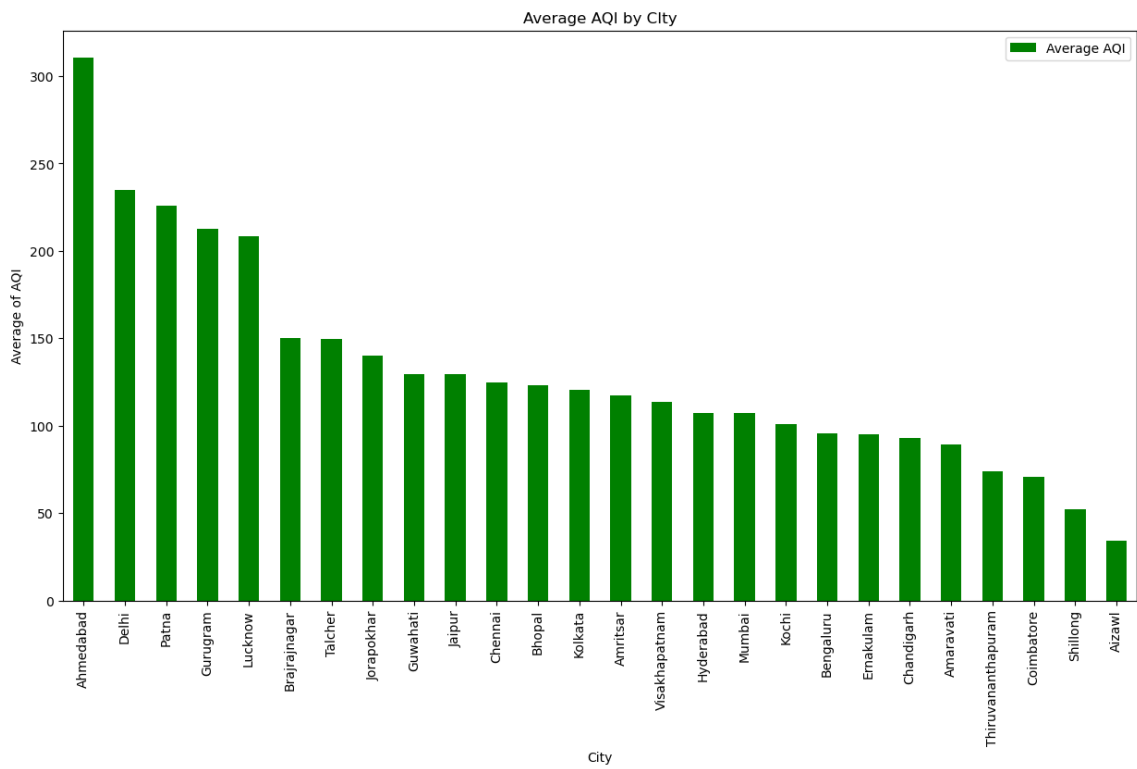


## Interpretations:

- From the heatmap, we can say that other pollutants do not affect the AQI much.
- They have the positive correlation with AQI but not strong correlation.
- PM2.5 and PM10 has good positive correlation between them.
- Also Toluene and Benzene have high correlation.

In [39]:

```
1 # Check the citywise average of AQI :
2
3 AQI_df.groupby('City')['AQI'].mean().sort_values(ascending=False).plot(kind='bar', color='green',
4 plt.ylabel('Average of AQI')
5 plt.title("Average AQI by City")
6 plt.legend()
7 plt.show()
```



In [57]:

```
1 dh=pd.read_csv("C:\\Users\\snaga\\OneDrive\\Desktop\\state.csv")
```

Out[57]:

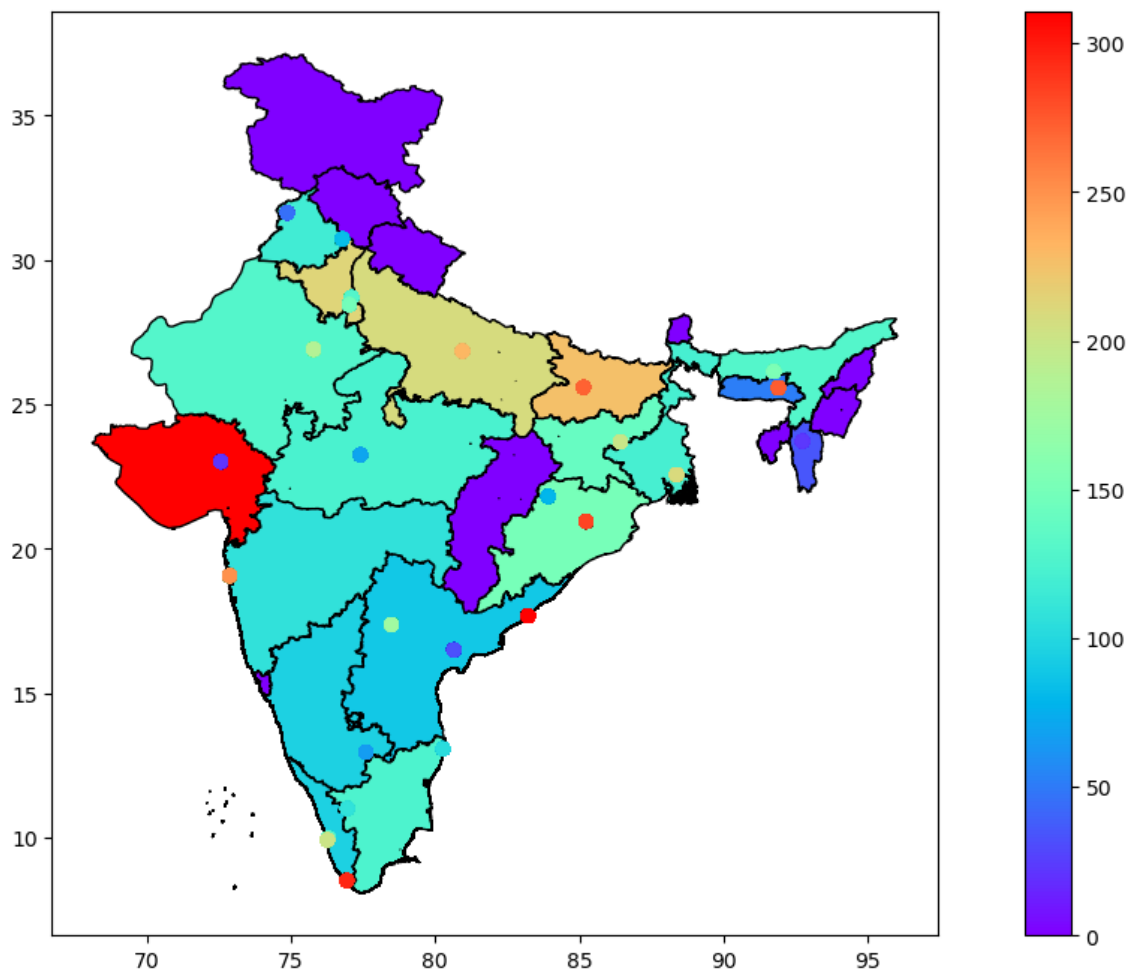
Unnamed: 0		ST_NAME	AQI_DATA
0	0	ANDAMAN AND NICOBAR ISLANDS	NaN
1	1	Andhra Pradesh	89.511415
2	2	Arunachal Pradesh	NaN
3	3	Assam	129.688454
4	4	Bihar	226.050314

In [58]:

```
1 #importing India shape file as df
2 df=gpd.read_file("C:\\Users\\snaga\\OneDrive\\Desktop\\in\\INDIA_states.shp")
3 df["AQI_VAL"]=dh["AQI_DATA"]
4
5 base=df.plot(column="AQI_VAL",cmap="rainbow",edgecolor='black',legend=True,)
6 geo.plot(ax=base,cmap="rainbow",)
7
```

Out[58]:

&lt;AxesSubplot:&gt;



In [39]:

```
1 # First seperate target and independent variables :
2 x = AQI_df.drop(['City', 'Date', 'AQI'], axis=1)
3 #x = sm.add_constant(x)
4 y= AQI_df.AQI
5
6 # Fitting the model
7 model = OLS(y, x).fit()
8 model.summary()
```

Out[39]:

## OLS Regression Results

<b>Dep. Variable:</b>	AQI	<b>R-squared (uncentered):</b>	0.972
<b>Model:</b>	OLS	<b>Adj. R-squared (uncentered):</b>	0.972
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	2733.
<b>Date:</b>	Thu, 23 Feb 2023	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	16:14:47	<b>Log-Likelihood:</b>	-4554.4
<b>No. Observations:</b>	949	<b>AIC:</b>	9133.
<b>Df Residuals:</b>	937	<b>BIC:</b>	9191.
<b>Df Model:</b>	12		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>PM2_5</b>	1.4955	0.027	55.651	0.000	1.443	1.548
<b>PM10</b>	0.1304	0.019	6.861	0.000	0.093	0.168
<b>NO</b>	0.7029	0.090	7.819	0.000	0.527	0.879
<b>NO2</b>	0.2631	0.079	3.335	0.001	0.108	0.418
<b>NOx</b>	-0.2420	0.051	-4.702	0.000	-0.343	-0.141
<b>NH3</b>	0.0794	0.073	1.094	0.274	-0.063	0.222
<b>CO</b>	13.0731	0.343	38.117	0.000	12.400	13.746
<b>SO2</b>	-0.2595	0.091	-2.859	0.004	-0.438	-0.081
<b>O3</b>	0.4186	0.057	7.291	0.000	0.306	0.531
<b>Benzene</b>	-0.4713	0.330	-1.429	0.153	-1.119	0.176
<b>Toluene</b>	-0.0012	0.146	-0.008	0.993	-0.288	0.286
<b>Xylene</b>	0.2218	0.250	0.887	0.375	-0.269	0.712

<b>Omnibus:</b>	112.188	<b>Durbin-Watson:</b>	0.781
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	935.496
<b>Skew:</b>	-0.133	<b>Prob(JB):</b>	7.24e-204
<b>Kurtosis:</b>	7.857	<b>Cond. No.</b>	52.8

## Notes:

- [1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correctly specified.

In [40]:

```

1 # Print significant variables which are most affect on AQI :
2
3 significant = model.pvalues[model.pvalues < 0.05].sort_values(ascending=True)
4
5 print('Significant vaeiables which more affect on AQI :\n', significant)

```

Significant vaeiables which more affect on AQI :

```

PM2_5      2.774664e-299
CO         1.034701e-192
NO         1.434162e-14
O3         6.541010e-13
PM10       1.238154e-11
NOx        2.966255e-06
NO2        8.860833e-04
SO2        4.341676e-03
dtype: float64

```

## Interpretations:

- The pollutants CO, PM2.5, NO2 and SO2 are affecting the most on AQI
- Benzene and NO are least significant.

## Multi Linear Regression Model

In [41]:

```

1 x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=3,test_size=0.2)
2 model=LinearRegression()
3 model.fit(x_train,y_train)
4
5 y_pred=model.predict(x_test)
6
7

```

In [42]:

```

1
2 print("coef of determination(accuracy of Model):",model.score(x_train,y_train)*100,"
3 print("*****")
4 print("intercept is:",model.intercept_)
5 print("*****")
6 mae1 = mean_absolute_error(y_test, y_pred)
7 print("mean absolute error",mae1)
8
9
10

```

coef of determination(accuracy of Model): 90.06172778136546 percent

\*\*\*\*\*

intercept is: 17.25179018772488

\*\*\*\*\*

mean absolute error 20.09092741263187



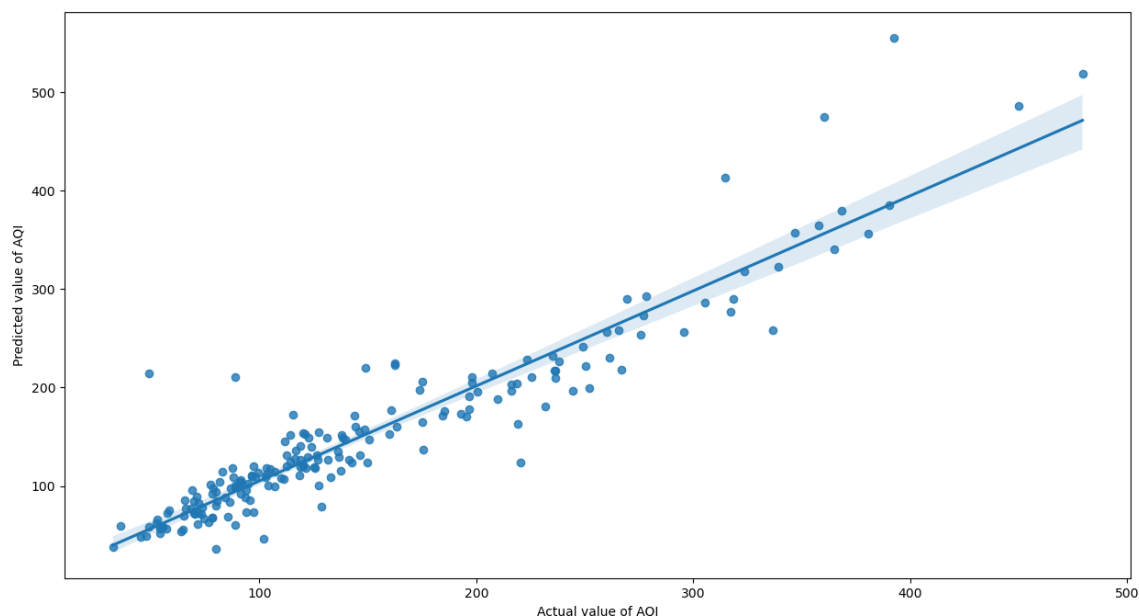
# Visualization

In [43]:

```
1 import matplotlib.pyplot as plt
2 sns.regplot(y_test,y_pred)
3 plt.xlabel('Actual value of AQI')
4 plt.ylabel('Predicted value of AQI')
5
```

Out[43]:

Text(0, 0.5, 'Predicted value of AQI')



## Random forest Regression

In [44]:

```
1 model = RandomForestRegressor(n_estimators=100, max_depth=5, random_state=42)
2 model.fit(x_train, y_train)
3 y_pred1 = model.predict(x_test)
4
5 r2 = r2_score(y_test, y_pred1)
6
7
8 print("Accuracy is:: ", r2*100,"percent")
9
10 print("*****")
11 mae2 = mean_absolute_error(y_test, y_pred1)
12 print("mean absolute error",mae2)
```

Accuracy is:: 93.57444362392091 percent

\*\*\*\*\*

mean absolute error 15.459628917216765

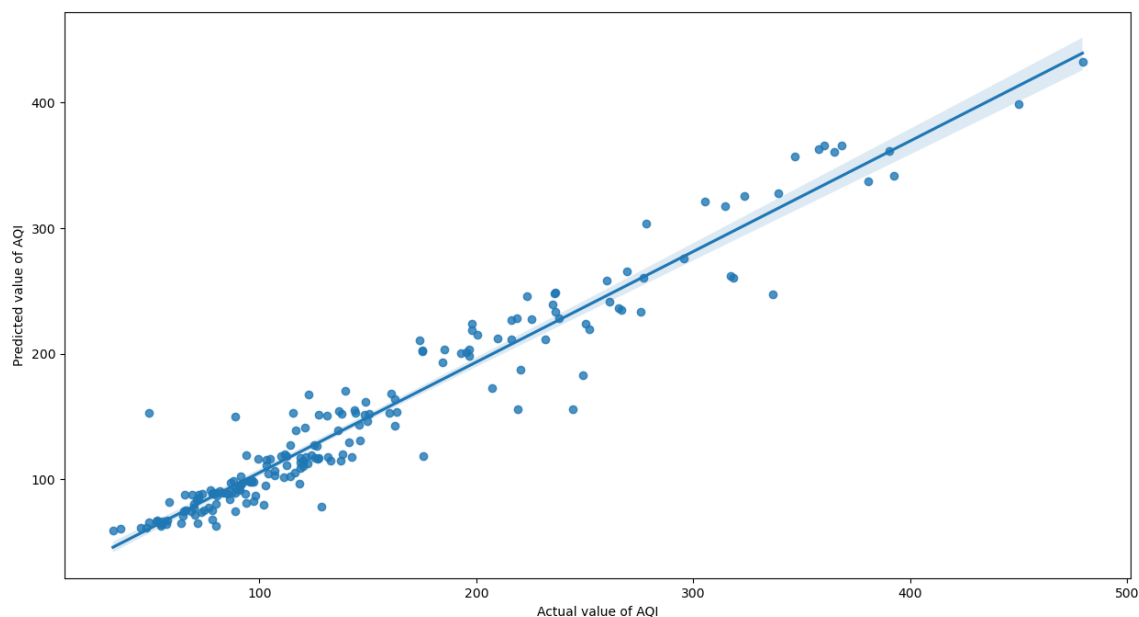
# Visualization

In [45]:

```
1 sns.regplot(y_test,y_pred1)
2 plt.xlabel('Actual value of AQI')
3 plt.ylabel('Predicted value of AQI')
4
```

Out[45]:

Text(0, 0.5, 'Predicted value of AQI')



## Decision Tree Regressor

In [46]:

```
1 # Create a decision tree regressor object
2 regressor = DecisionTreeRegressor(max_depth=3,random_state=20)
3
4 # Fit the regressor with the training data
5 regressor.fit(x_train, y_train)
6
7 # Make predictions on the testing data
8 y_pred2 = regressor.predict(x_test)
9 r4 = r2_score(y_test, y_pred2)
10 print("Accuracy is:",r4*100,"percent")
11 print("*****")
12
13 mae3 = mean_absolute_error(y_test, y_pred2)
14 print("mean absolute error",mae3)
```

Accuracy is: 82.24528017778303 percent

\*\*\*\*\*

mean absolute error 26.20649765527541

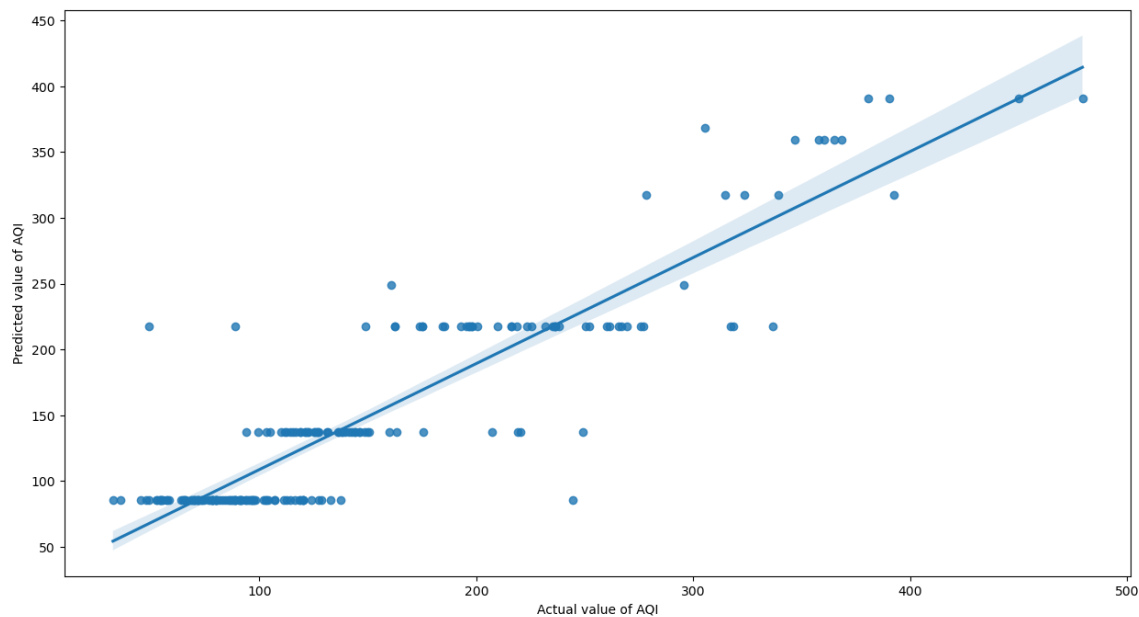
# Visualization

In [47]:

```
1 sns.regplot(y_test,y_pred2)
2 plt.xlabel('Actual value of AQI')
3 plt.ylabel('Predicted value of AQI')
4
```

Out[47]:

Text(0, 0.5, 'Predicted value of AQI')



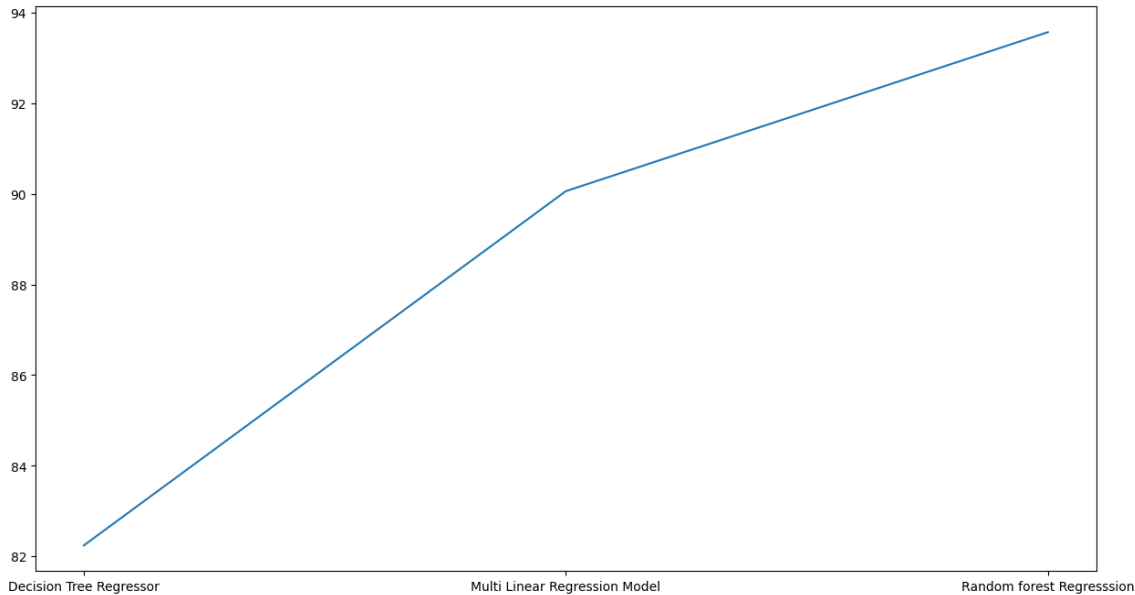
## Comparison

In [48]:

```
1 a=["Decision Tree Regressor","Multi Linear Regression Model","Random forest Regressor"]
2 b=[82.24,90.06,93.57]
3 sns.lineplot(a,b)
```

Out[48]:

&lt;AxesSubplot:&gt;



## CONCLUSIONS

- We have done a Multivariate analysis on AQI.
- We have recieved the best results for predicting AQI through random forest regression, therefore we chose Choose as our final model and the forecasted values are visualized through graphical representation.
- From the forecasted graph we can conclude that AQI will be high for coming years. The focus should be on how the pollutants affecting AQI can be reduced.
- According to our analysis the pollutants affecting the most on AQI are PM2.5, NO2 and CO on the basis of Correlation and Statistical Analysis.

In [ ]:

1