

Importing Libraries

In [1]:

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

Importing Datasets

In [2]:

```
df=pd.read_csv(r"C:\Users\user\Downloads\C10_air\csvs_per_year\csvs(Dataset)\madrid_2006.
df
```

Out[2]:

	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3	
0	2006-02-01 01:00:00	NaN	1.84	NaN	NaN	NaN	155.100006	490.100006	NaN	4.880000	97
1	2006-02-01 01:00:00	1.68	1.01	2.38	6.36	0.32	94.339996	229.699997	3.04	7.100000	25
2	2006-02-01 01:00:00	NaN	1.25	NaN	NaN	NaN	66.800003	192.000000	NaN	4.430000	34
3	2006-02-01 01:00:00	NaN	1.68	NaN	NaN	NaN	103.000000	407.799988	NaN	4.830000	28
4	2006-02-01 01:00:00	NaN	1.31	NaN	NaN	NaN	105.400002	269.200012	NaN	6.990000	54
...
230563	2006-05-01 00:00:00	5.88	0.83	6.23	NaN	0.20	112.500000	218.000000	NaN	24.389999	93
230564	2006-05-01 00:00:00	0.76	0.32	0.48	1.09	0.08	51.900002	54.820000	0.61	48.410000	29
230565	2006-05-01 00:00:00	0.96	NaN	0.69	NaN	0.19	135.100006	179.199997	NaN	11.460000	64
230566	2006-05-01 00:00:00	0.50	NaN	0.67	NaN	0.10	82.599998	105.599998	NaN	NaN	94
230567	2006-05-01 00:00:00	1.95	0.74	1.99	4.00	0.24	107.300003	160.199997	2.01	17.730000	52

230568 rows × 17 columns

Data Cleaning and Data Preprocessing

In [3]:

```
df=df.dropna()
```

In [4]:

```
df.columns
```

Out[4]:

```
Index(['date', 'BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
      'PM10', 'PM25', 'PXY', 'SO_2', 'TCH', 'TOL', 'station'],
      dtype='object')
```

In [5]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 24758 entries, 5 to 230567
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        24758 non-null  object
1   BEN         24758 non-null  float64
2   CO          24758 non-null  float64
3   EBE         24758 non-null  float64
4   MXY         24758 non-null  float64
5   NMHC        24758 non-null  float64
6   NO_2        24758 non-null  float64
7   NOx         24758 non-null  float64
8   OXY         24758 non-null  float64
9   O_3         24758 non-null  float64
10  PM10        24758 non-null  float64
11  PM25        24758 non-null  float64
12  PXY         24758 non-null  float64
13  SO_2        24758 non-null  float64
14  TCH         24758 non-null  float64
15  TOL         24758 non-null  float64
16  station     24758 non-null  int64
dtypes: float64(15), int64(1), object(1)
memory usage: 3.4+ MB
```

In [7]:

```
data=df[['NMHC', 'NO_2', 'O_3']]
data
```

Out[7]:

	NMHC	NO_2	O_3
5	0.44	142.199997	5.990000
22	0.17	59.910000	2.450000
25	0.40	117.699997	4.780000
31	0.25	92.059998	5.920000
48	0.16	60.189999	2.280000
...
230538	0.10	49.259998	64.599998
230541	0.33	63.220001	17.670000
230547	0.26	202.399994	11.130000
230564	0.08	51.900002	48.410000
230567	0.24	107.300003	17.730000

24758 rows × 3 columns

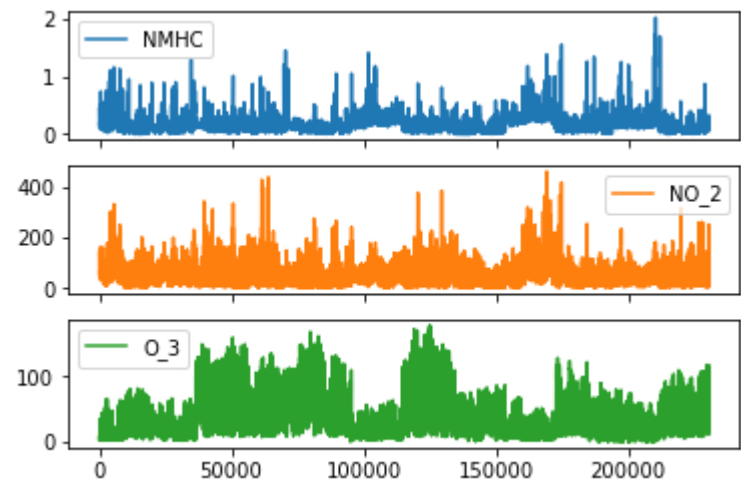
Line chart

In [8]:

```
data.plot.line(subplots=True)
```

Out[8]:

array([<AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>], dtype=object)



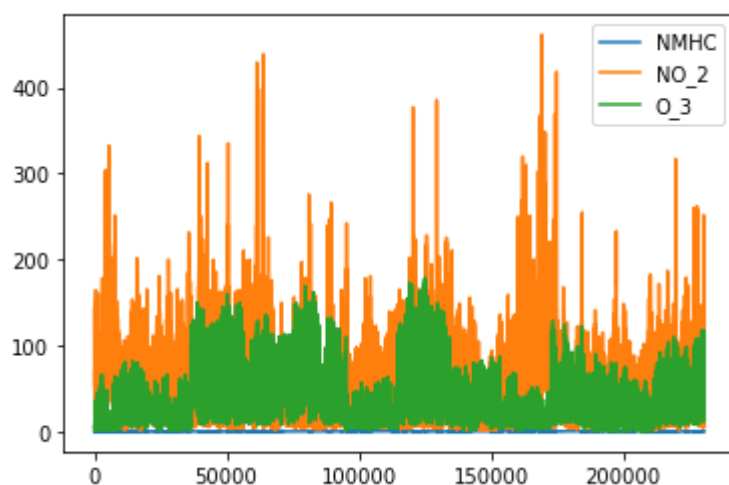
Line chart

In [10]:

```
data.plot.line()
```

Out[10]:

<AxesSubplot:>



Bar chart

In [11]:

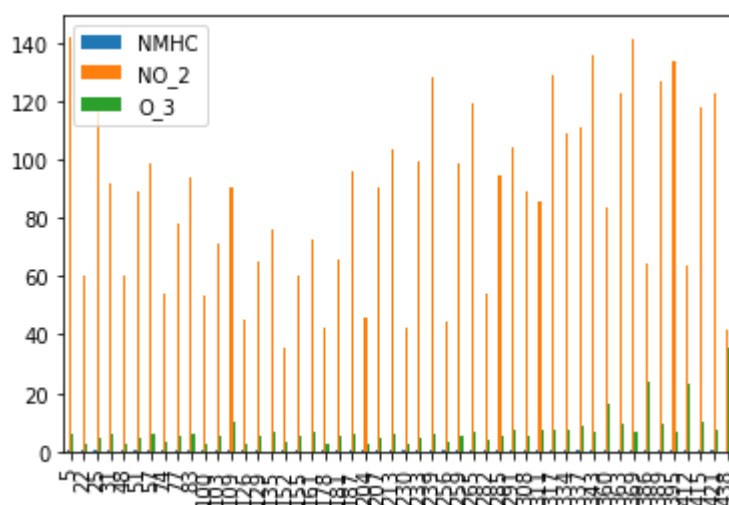
```
b=data[0:50]
```

In [12]:

```
b.plot.bar()
```

Out[12]:

<AxesSubplot:>



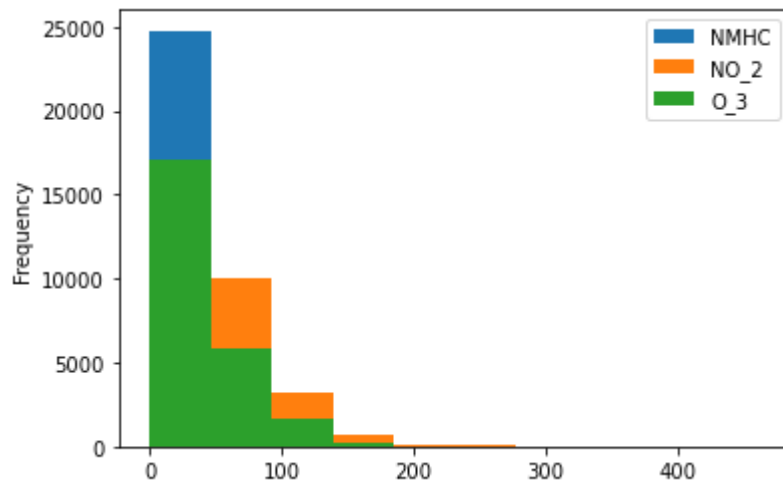
Histogram

In [13]:

```
data.plot.hist()
```

Out[13]:

<AxesSubplot:ylabel='Frequency'>



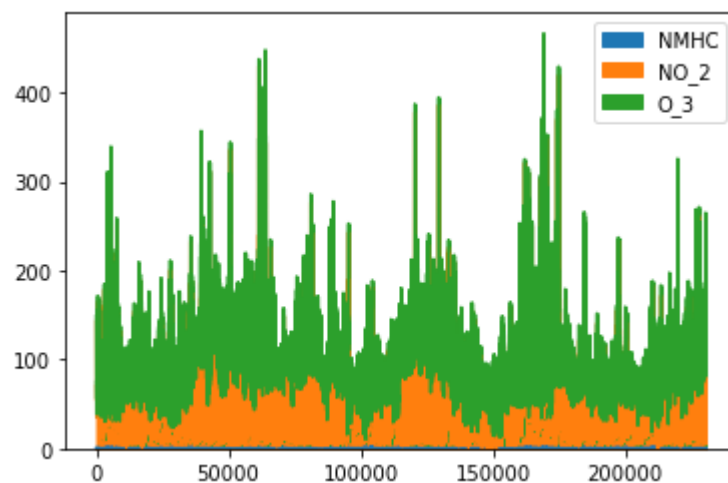
Area chart

In [14]:

```
data.plot.area()
```

Out[14]:

<AxesSubplot:>



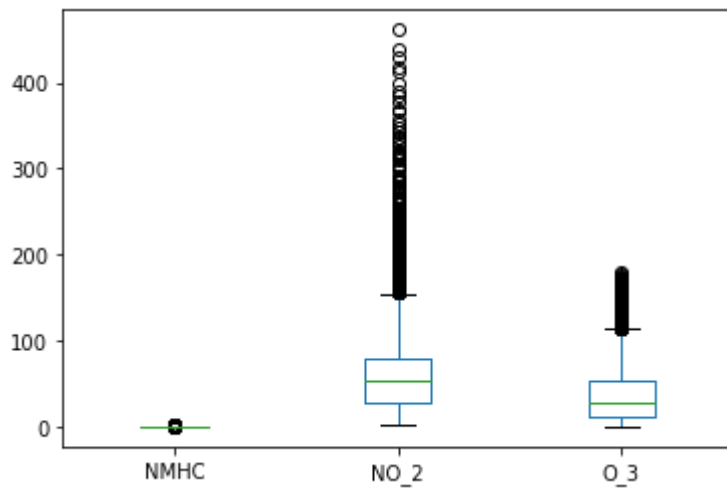
Box chart

In [15]:

```
data.plot.box()
```

Out[15]:

<AxesSubplot:>



Pie chart

In [17]:

```
b.plot.pie(y='O_3' )
```

Out[17]:

<AxesSubplot:ylabel='O_3'>



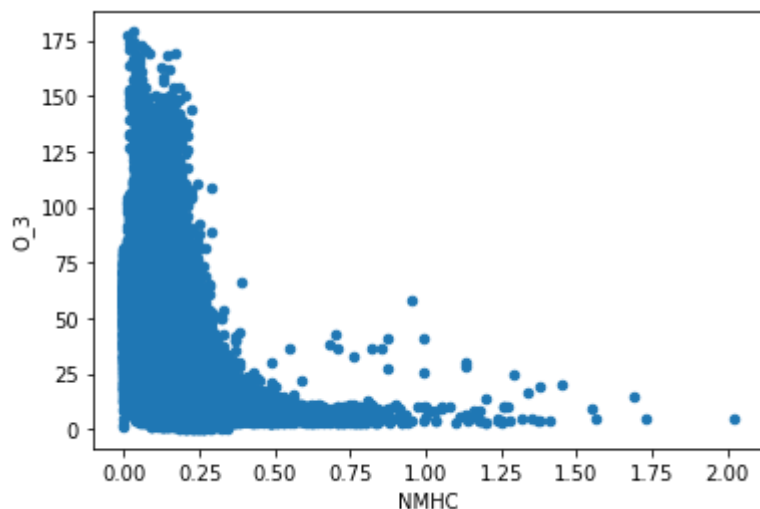
Scatter chart

In [18]:

```
data.plot.scatter(x='NMHC', y='O_3')
```

Out[18]:

```
<AxesSubplot:xlabel='NMHC', ylabel='O_3'>
```



In [19]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 24758 entries, 5 to 230567
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        24758 non-null  object
1   BEN         24758 non-null  float64
2   CO          24758 non-null  float64
3   EBE         24758 non-null  float64
4   MXY         24758 non-null  float64
5   NMHC        24758 non-null  float64
6   NO_2        24758 non-null  float64
7   NOx         24758 non-null  float64
8   OXY         24758 non-null  float64
9   O_3         24758 non-null  float64
10  PM10        24758 non-null  float64
11  PM25        24758 non-null  float64
12  PXY         24758 non-null  float64
13  SO_2        24758 non-null  float64
14  TCH         24758 non-null  float64
15  TOL         24758 non-null  float64
16  station     24758 non-null  int64
dtypes: float64(15), int64(1), object(1)
memory usage: 3.4+ MB
```

In [20]:

```
df.describe()
```

Out[20]:

	BEN	CO	EBE	MXY	NMHC	NO_2
count	24758.000000	24758.000000	24758.000000	24758.000000	24758.000000	24758.000000
mean	1.350624	0.600713	1.824534	3.835034	0.176546	58.333481
std	1.541636	0.419048	1.868939	4.069036	0.126683	40.529382
min	0.110000	0.000000	0.170000	0.150000	0.000000	1.680000
25%	0.450000	0.360000	0.810000	1.060000	0.100000	28.450001
50%	0.850000	0.500000	1.130000	2.500000	0.150000	52.959999
75%	1.680000	0.720000	2.160000	5.090000	0.220000	79.347498
max	45.430000	7.250000	57.799999	66.900002	2.020000	461.299988

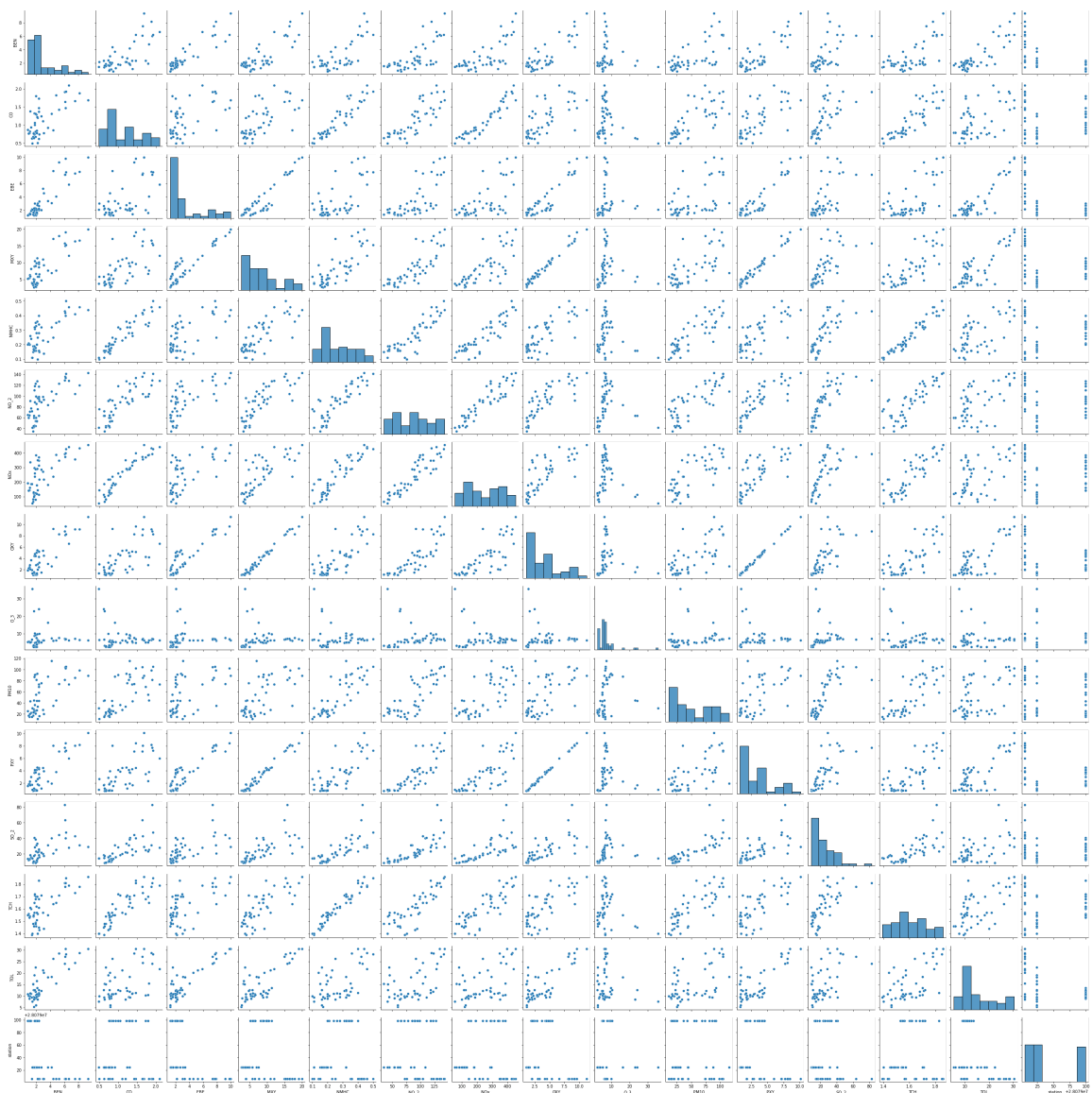
In [21]:

```
df1=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',  
        'PM10', 'PXY', 'SO_2', 'TCH', 'TOL', 'station']]
```

EDA AND VISUALIZATION

```
sns.pairplot(df1[0:50])
```

```
<seaborn.axisgrid.PairGrid at 0x21a7763ad30>
```



In [67]:

```
x=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',  
      'PM10', 'PXY', 'SO_2', 'TCH', 'TOL']]  
y=df['station']
```

In [68]:

```
from sklearn.model_selection import train_test_split  
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

Linear Regression

In [69]:

```
from sklearn.linear_model import LinearRegression  
lr=LinearRegression()  
lr.fit(x_train,y_train)
```

Out[69]:

LinearRegression()

In [70]:

```
lr.intercept_
```

Out[70]:

28079021.191809315

In [71]:

```
coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

Out[71]:

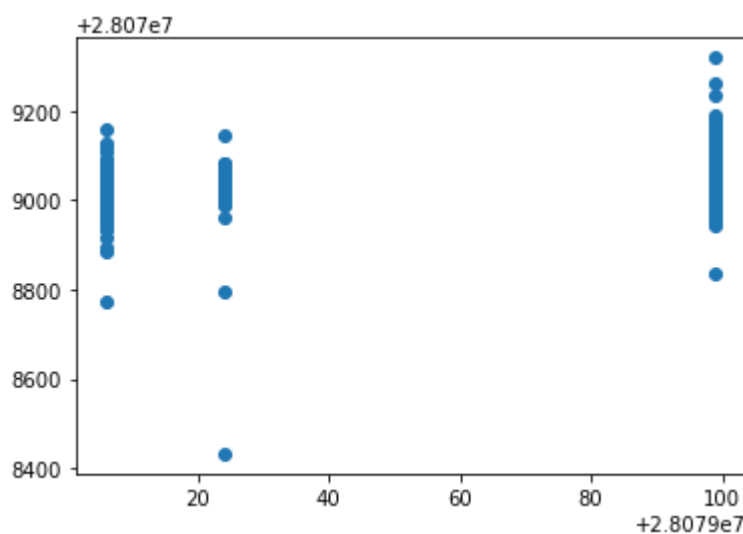
	Co-efficient
BEN	-18.378990
CO	-9.956654
EBE	-23.531143
MXY	4.689670
NMHC	126.649643
NO_2	-0.016238
NOx	-0.004541
OXY	15.492118
O_3	-0.053784
PM10	0.131450
PXY	5.674572
SO_2	-0.652901
TCH	17.650518
TOL	-0.517190

In [72]:

```
prediction =lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[72]:

<matplotlib.collections.PathCollection at 0x21a07f864f0>



ACCURACY

In [73]:

```
lr.score(x_test,y_test)
```

Out[73]:

0.38347413896089644

In [74]:

```
lr.score(x_train,y_train)
```

Out[74]:

0.3972450971871986

Ridge and Lasso

In [75]:

```
from sklearn.linear_model import Ridge,Lasso
```

In [76]:

```
rr=Ridge(alpha=10)  
rr.fit(x_train,y_train)
```

Out[76]:

Ridge(alpha=10)

Accuracy(Ridge)

In [77]:

```
rr.score(x_test,y_test)
```

Out[77]:

0.3820493845006965

In [78]:

```
rr.score(x_train,y_train)
```

Out[78]:

0.39661106884033615

In [79]:

```
la=Lasso(alpha=10)  
la.fit(x_train,y_train)
```

Out[79]:

Lasso(alpha=10)

In [81]:

```
la.score(x_test,y_test)
```

Out[81]:

0.05647819697010681

Accuracy(Lasso)

In [80]:

```
la.score(x_train,y_train)
```

Out[80]:

0.062336248080204326

Accuracy(Elastic Net)

In [82]:

```
from sklearn.linear_model import ElasticNet  
en=ElasticNet()  
en.fit(x_train,y_train)
```

Out[82]:

ElasticNet()

In [83]:

```
en.coef_
```

Out[83]:

```
array([-8.55559241e+00,  0.00000000e+00, -9.05278481e+00,  3.41965661e+00,  
        4.01206179e-01, -3.00597081e-03,  5.10526761e-03,  3.44982093e+00,  
       -1.23661221e-01,  2.93481571e-01,  2.43724729e+00, -4.22212620e-01,  
        5.28019988e-01, -9.96915988e-01])
```

In [84]:

```
en.intercept_
```

Out[84]:

28079052.093522523

In [85]:

```
prediction=en.predict(x_test)
```


In [86]:

```
en.score(x_test,y_test)
```

Out[86]:

0.2328938395355542

Evaluation Metrics

In [87]:

```
from sklearn import metrics
print(metrics.mean_absolute_error(y_test,prediction))
print(metrics.mean_squared_error(y_test,prediction))
print(np.sqrt(metrics.mean_squared_error(y_test,prediction)))
```

32.338665036423414

1266.5452114192653

35.58855450027811

Logistic Regression

In [88]:

```
from sklearn.linear_model import LogisticRegression
```

In [89]:

```
feature_matrix=df[['BEN', 'CO', 'EBE', 'MXV', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
                  'PM10', 'PXY', 'SO_2', 'TCH', 'TOL']]
target_vector=df[ 'station']
```

In [90]:

```
feature_matrix.shape
```

Out[90]:

(24758, 14)

In [91]:

```
target_vector.shape
```

Out[91]:

(24758,)

In [92]:

```
from sklearn.preprocessing import StandardScaler
```

In [93]:

```
fs=StandardScaler().fit_transform(feature_matrix)
```

In [94]:

```
logr=LogisticRegression(max_iter=10000)  
logr.fit(fs,target_vector)
```

Out[94]:

```
LogisticRegression(max_iter=10000)
```

In [95]:

```
observation=[[1,2,3,4,5,6,7,8,9,10,11,12,13,14]]
```

In [96]:

```
prediction=logr.predict(observation)  
print(prediction)
```

```
[28079099]
```

In [97]:

```
logr.classes_
```

Out[97]:

```
array([28079006, 28079024, 28079099], dtype=int64)
```

In [98]:

```
logr.score(fs,target_vector)
```

Out[98]:

```
0.8741416915744405
```

In [99]:

```
logr.predict_proba(observation)[0][0]
```

Out[99]:

```
3.5557727473608076e-15
```

In [100]:

```
logr.predict_proba(observation)
```

Out[100]:

```
array([[3.55577275e-15, 7.80743173e-29, 1.00000000e+00]])
```

Random Forest

In [101]:

```
from sklearn.ensemble import RandomForestClassifier
```

In [102]:

```
rfc=RandomForestClassifier()  
rfc.fit(x_train,y_train)
```

Out[102]:

```
RandomForestClassifier()
```

In [103]:

```
parameters={'max_depth':[1,2,3,4,5],  
            'min_samples_leaf':[5,10,15,20,25],  
            'n_estimators':[10,20,30,40,50]}  
}
```

In [104]:

```
from sklearn.model_selection import GridSearchCV  
grid_search =GridSearchCV(estimator=rfc,param_grid=parameters,cv=2,scoring="accuracy")  
grid_search.fit(x_train,y_train)
```

Out[104]:

```
GridSearchCV(cv=2, estimator=RandomForestClassifier(),  
             param_grid={'max_depth': [1, 2, 3, 4, 5],  
                         'min_samples_leaf': [5, 10, 15, 20, 25],  
                         'n_estimators': [10, 20, 30, 40, 50]}},  
             scoring='accuracy')
```

In [105]:

```
grid_search.best_score_
```

Out[105]:

```
0.8757068667051355
```

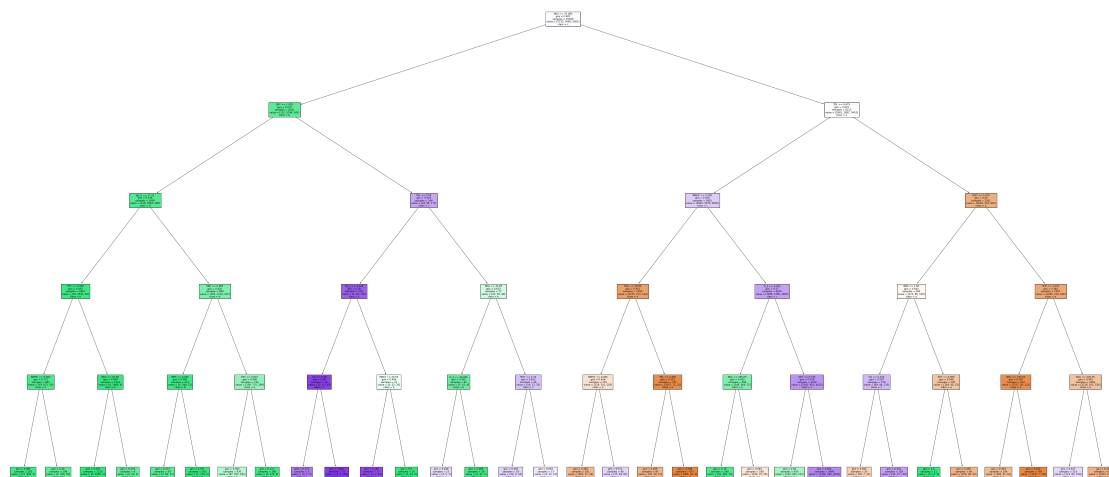
In [106]:

```
rfc_best=grid_search.best_estimator_
```

In [107]:

```
from sklearn.tree import plot_tree

plt.figure(figsize=(80,40))
plot_tree(rfc_best.estimators_[5],feature_names=x.columns,class_names=['a','b','c','d'],f
[113, 64, 159]\nclclass = 'c'),
Text(4394.25, 181.19999999999982, 'gini = 0.42\nsamples = 850\nvalue =
[1015, 37, 365]\nclclass = 'a')]
```



Conclusion

Accuracy

Linear Regression:0.3972450971871986

Ridge Regression:0.39661106884033615

Lasso Regression:0.062336248080204326

ElasticNet Regression:0.2328938395355542

Logistic Regression:0.8741416915744405

Random Forest:0.8757068667051355

Random Forest is suitable for this dataset