```
In [1]: import numpy as np
        import pandas as pd
        import seaborn as sns
        import matplotlib.pyplot as plt
```

```
In [2]: df=pd.read_csv("madrid_2011.csv")
```

```
In [3]: df.head()
```

Out[3]:

| | date | BEN | CO | EBE | NMHC | NO | NO_2 | O_3 | PM10 | PM25 | SO_2 | TCH | TOL | stat |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2011-11-01 01:00:00 | NaN | 1.0 | NaN | NaN | 154.0 | 84.0 | NaN | NaN | NaN | 6.0 | NaN | NaN | 28079 |
| 1 | 2011-11-01 01:00:00 | 2.5 | 0.4 | 3.5 | 0.26 | 68.0 | 92.0 | 3.0 | 40.0 | 24.0 | 9.0 | 1.54 | 8.7 | 28079 |
| 2 | 2011-11-01 01:00:00 | 2.9 | NaN | 3.8 | NaN | 96.0 | 99.0 | NaN | NaN | NaN | NaN | NaN | 7.2 | 28079 |
| 3 | 2011-11-01 01:00:00 | NaN | 0.6 | NaN | NaN | 60.0 | 83.0 | 2.0 | NaN | NaN | NaN | NaN | NaN | 28079 |
| 4 | 2011-11-01 01:00:00 | NaN | NaN | NaN | NaN | 44.0 | 62.0 | 3.0 | NaN | NaN | 3.0 | NaN | NaN | 28079 |

```
In [4]: df=df.dropna()
```

```
In [5]: df.columns
```

Out[5]: Index(['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',
               'SO_2', 'TCH', 'TOL', 'station'],
              dtype='object')

In [6]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 16460 entries, 1 to 209910
Data columns (total 14 columns):
 #   Column   Non-Null Count  Dtype
---  ------   --------------  -----
 0   date     16460 non-null  object
 1   BEN      16460 non-null  float64
 2   CO       16460 non-null  float64
 3   EBE      16460 non-null  float64
 4   NMHC     16460 non-null  float64
 5   NO       16460 non-null  float64
 6   NO_2     16460 non-null  float64
 7   O_3      16460 non-null  float64
 8   PM10     16460 non-null  float64
 9   PM25     16460 non-null  float64
 10  SO_2     16460 non-null  float64
 11  TCH      16460 non-null  float64
 12  TOL      16460 non-null  float64
 13  station  16460 non-null  int64
dtypes: float64(12), int64(1), object(1)
memory usage: 1.9+ MB
```
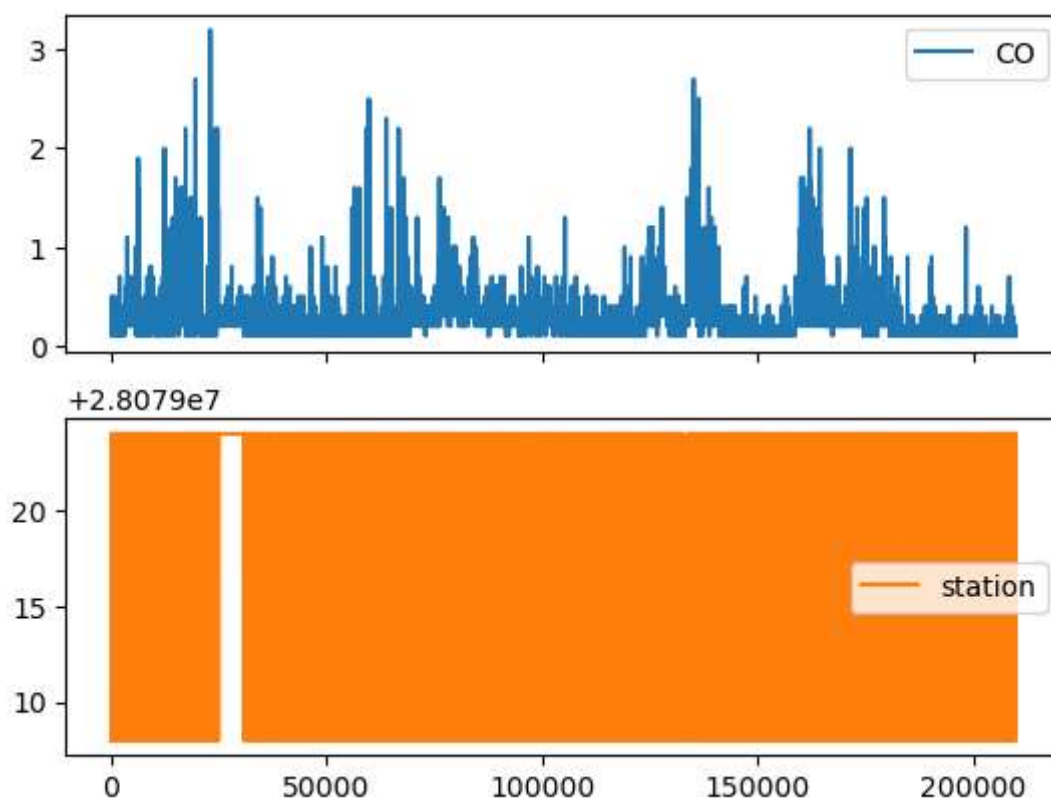
In [7]: `data=df[['CO','station']]`
`data`

Out[7]:

|        | CO  | station  |
|--------|-----|----------|
| 1      | 0.4 | 28079008 |
| 6      | 0.3 | 28079024 |
| 25     | 0.3 | 28079008 |
| 30     | 0.4 | 28079024 |
| 49     | 0.2 | 28079008 |
| ...    | ... | ...      |
| 209862 | 0.1 | 28079024 |
| 209881 | 0.1 | 28079008 |
| 209886 | 0.1 | 28079024 |
| 209905 | 0.1 | 28079008 |
| 209910 | 0.1 | 28079024 |

16460 rows × 2 columns
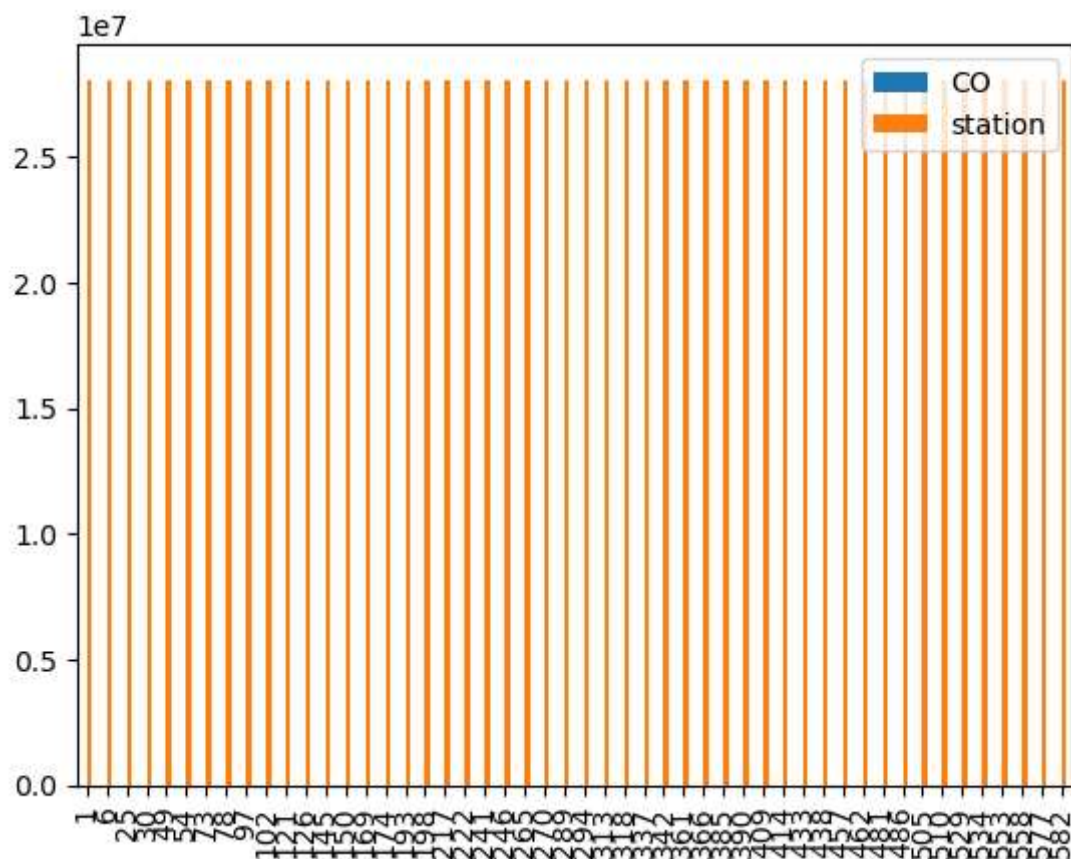
In [8]: `data.plot.line(subplots=True)`
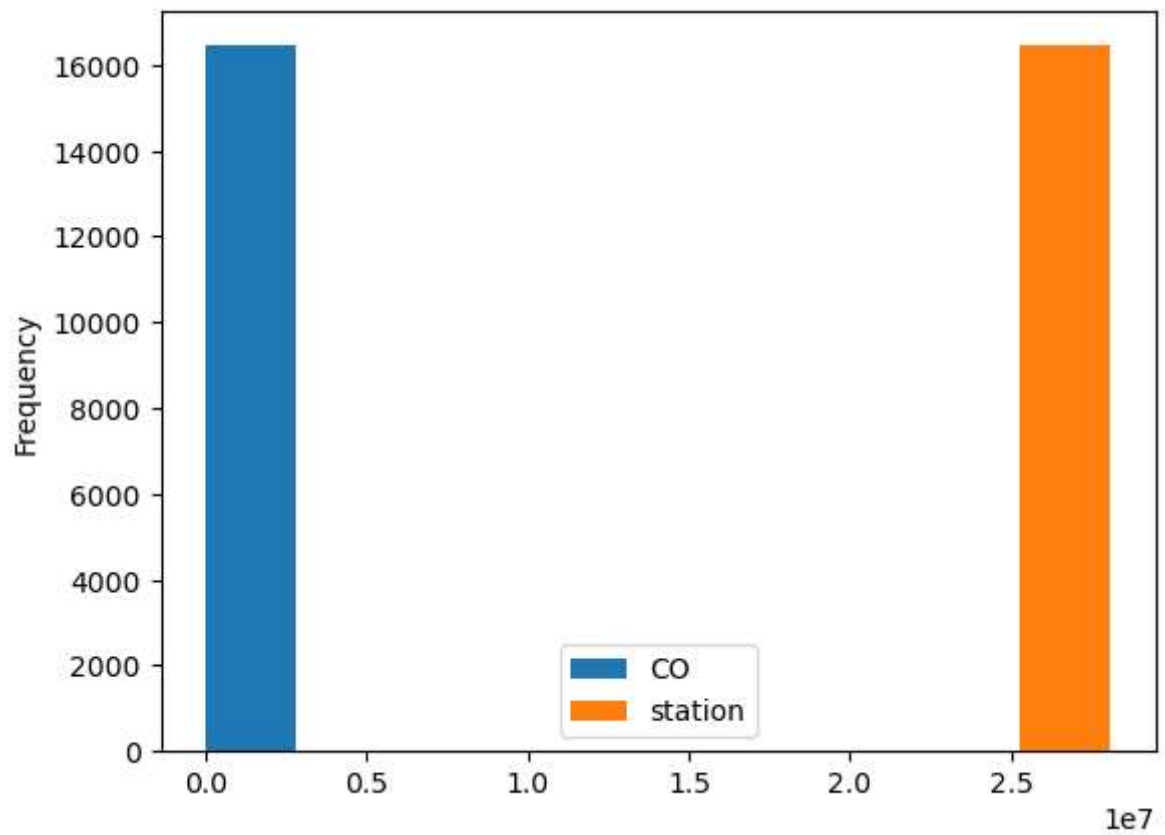
Out[8]: array([<Axes: >, <Axes: >], dtype=object)

In [9]:
```python
b=data[0:50]
b.plot.bar()
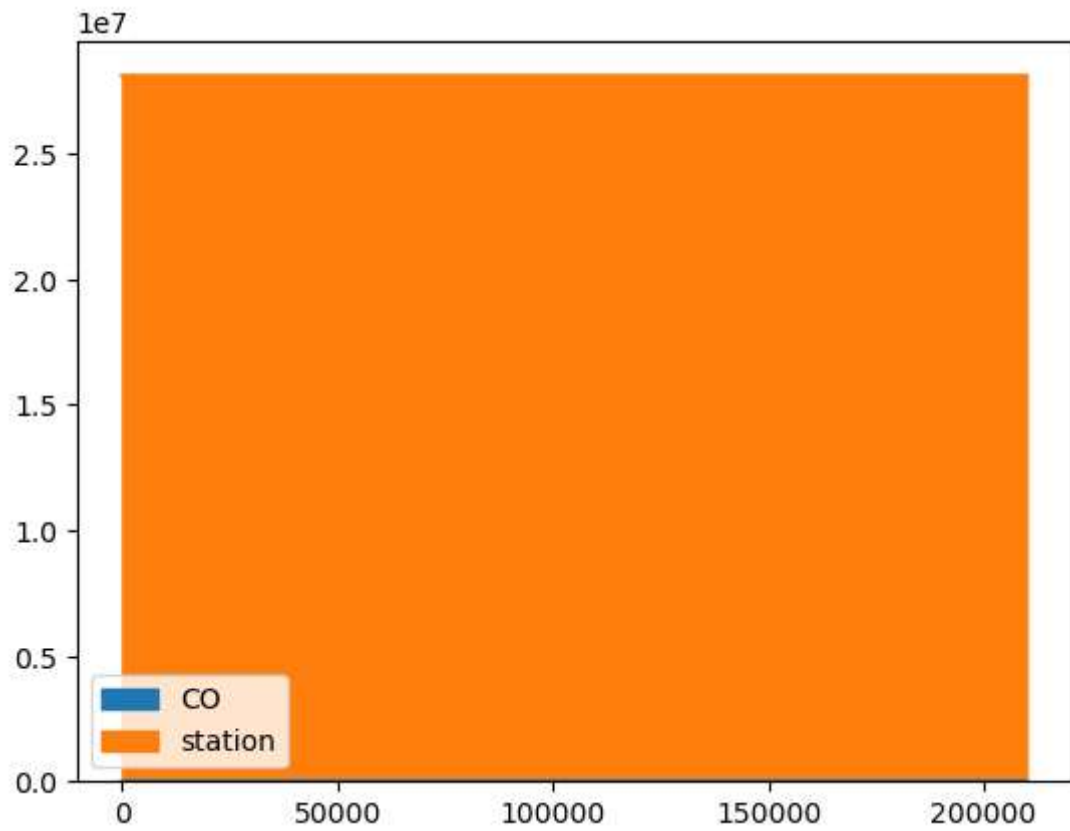```

Out[9]: <Axes: >

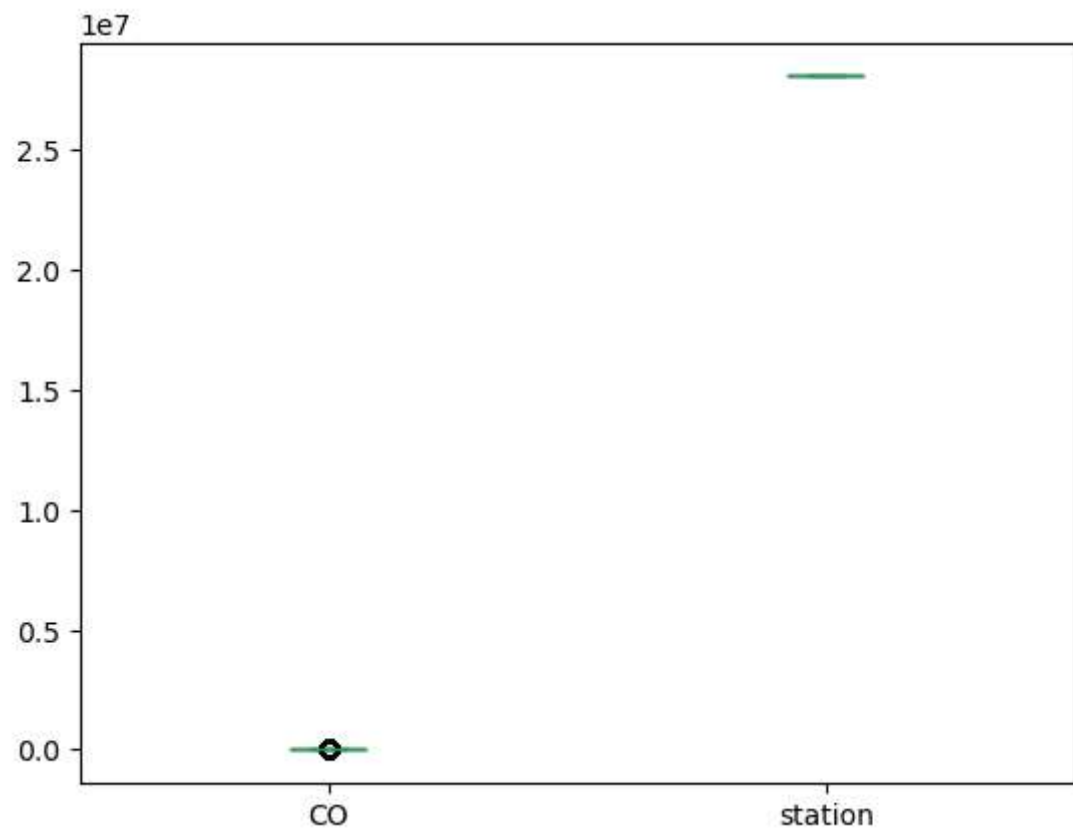In [10]: `data.plot.hist()`

Out[10]: <Axes: ylabel='Frequency'>

In [11]: `data.plot.area()`

Out[11]: <Axes: >
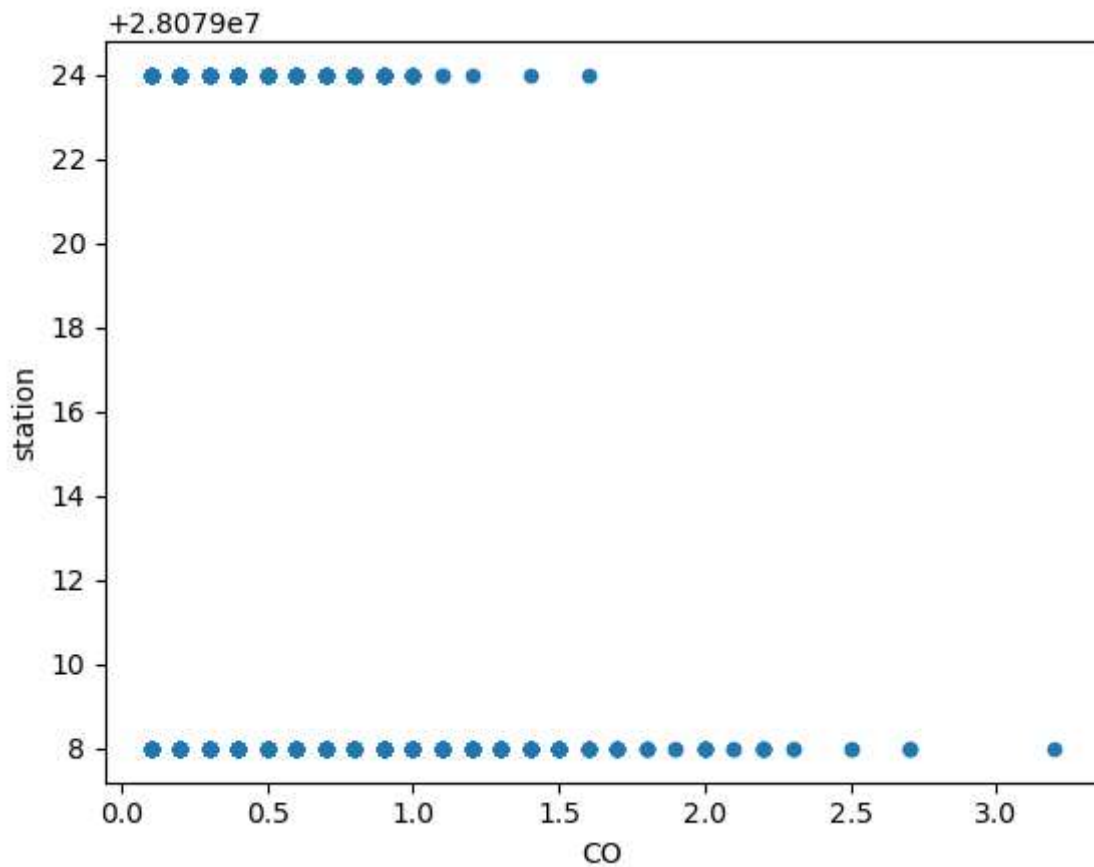
In [12]: `data.plot.box()`

Out[12]: `<Axes: >`

```
In [13]: data.plot.scatter(x='CO',y='station')
```

```
Out[13]: <Axes: xlabel='CO', ylabel='station'>
```



```
In [14]: x=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'NO', 'O_3',
         'PM10','PM25','SO_2', 'TCH', 'TOL']]
         y=df['station']
```
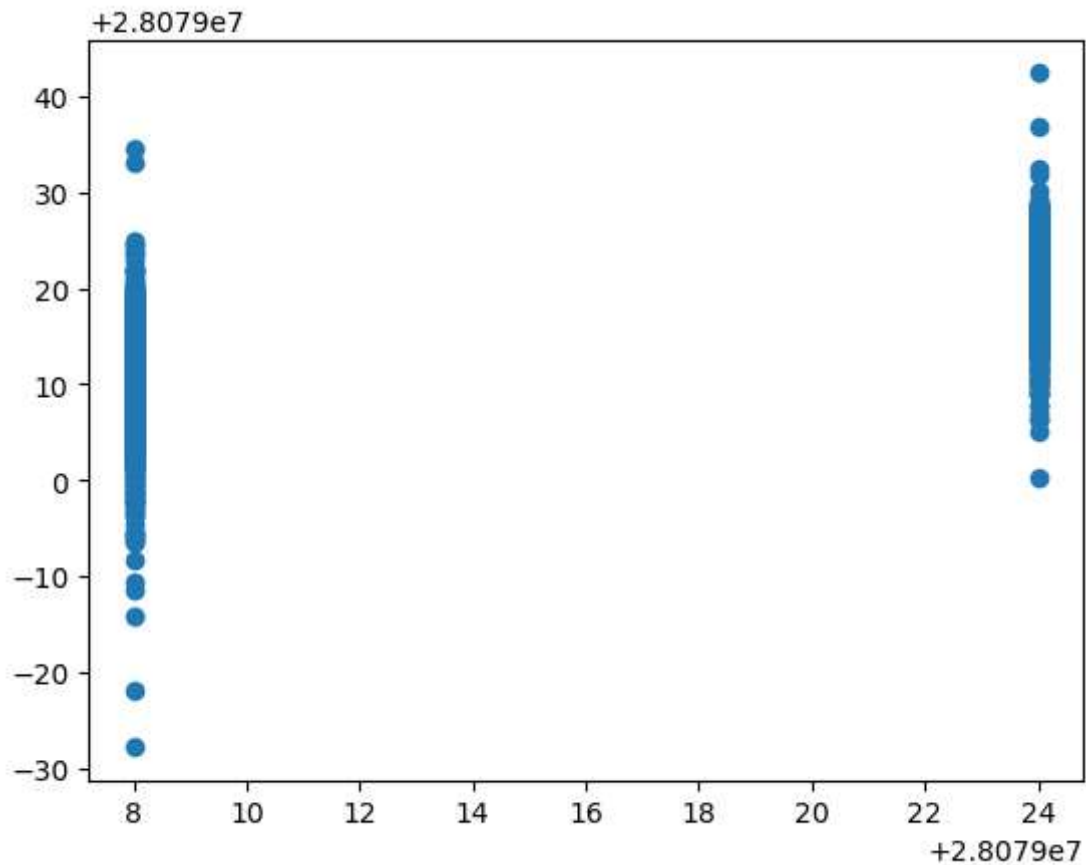
```
In [15]: from sklearn.model_selection import train_test_split
         x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

# Linear Regression

In [16]:
```python
from sklearn.linear_model import LinearRegression
lr=LinearRegression()
lr.fit(x_train,y_train)
lr.intercept_
prediction =lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[16]: <matplotlib.collections.PathCollection at 0x21c55761190>



In [17]:
```python
print(lr.score(x_test,y_test))
print(lr.score(x_train,y_train))
```

0.6259709949284429
0.6276479573022269

# Ridge and Lasso

```
In [18]: from sklearn.linear_model import Ridge,Lasso
         rr=Ridge(alpha=10)
         rr.fit(x_train,y_train)
         print(rr.score(x_test,y_test))
         print(rr.score(x_train,y_train))
         la=Lasso(alpha=10)
         la.fit(x_train,y_train)
```

```
0.5877091721395884
0.594469525578609
```

Out[18]:
```
 ▼      Lasso

Lasso(alpha=10)
```

```
In [19]: la.score(x_test,y_test)
```

Out[19]:  0.2273508727383743

# ElasticNet

```
In [20]: from sklearn.linear_model import ElasticNet
         en=ElasticNet()
         en.fit(x_train,y_train)
```

Out[20]:
```
 ▼ ElasticNet

ElasticNet()
```

```
In [21]: en.coef_
```

Out[21]:  array([ 0.2653043 ,  0.        , -0.        , -0.        , -0.14076021,
                0.05067364, -0.04412062,  0.02581064,  0.10485685, -0.16710751,
                0.        , -0.95067196])

```
In [22]: en.intercept_
```

Out[22]:  28079025.14184462

```
In [23]: prediction=en.predict(x_test)
```

```
In [24]: en.score(x_test,y_test)
```

Out[24]:  0.3356452219334174

# Evaluation Metrics

```
In [25]:  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)))
```

```
5.701842582725372
42.51283990709718
6.520187106755233
```

# Logistics Regression

```
In [26]:  from sklearn.linear_model import LogisticRegression
```

```
In [27]:  feature_matrix=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'NO', 'O_3',
          'PM10','PM25','SO_2', 'TCH', 'TOL']]
          target_vector=df[ 'station']
```

```
In [28]:  from sklearn.preprocessing import StandardScaler
          fs=StandardScaler().fit_transform(feature_matrix)
          logr=LogisticRegression(max_iter=10000)
          logr.fit(fs,target_vector)
```

```
Out[28]:  ▼        LogisticRegression
          LogisticRegression(max_iter=10000)
```

```
In [30]:  observation=[[1,2,3,4,5,6,7,8,9,10,11,12]]
          logr.predict_proba(observation)
```

```
Out[30]:  array([[1.00000000e+00, 1.13637857e-16]])
```

# Random Forest

```
In [31]:  from sklearn.ensemble import RandomForestClassifier
          rfc=RandomForestClassifier()
          rfc.fit(x_train,y_train)
```

```
Out[31]:  ▼ RandomForestClassifier
          RandomForestClassifier()
```

In [32]:
```python
parameters={'max_depth':[1,2,3,4,5],
'min_samples_leaf':[5,10,15,20,25],
'n_estimators':[10,20,30,40,50]
}
```

In [33]:
```python
from sklearn.model_selection import GridSearchCV
grid_search =GridSearchCV(estimator=rfc,param_grid=parameters,cv=2,scoring="ac
grid_search.fit(x_train,y_train)
```

Out[33]:

>        **GridSearchCV**
>
> ▸ **estimator: RandomForestClassifier**
>
>      ▸ RandomForestClassifier

In [34]:
```python
rfc_best=grid_search.best_estimator_
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'
```

```
\nvalue = [81, 8]\nclass = a'),
 Text(0.581395348372093, 0.25, 'TOL <= 1.55\ngini = 0.017\nsamples = 769
\nvalue = [1239, 11]\nclass = a'),
 Text(0.5581395348837209, 0.08333333333333333, 'gini = 0.107\nsamples = 12
1\nvalue = [182, 11]\nclass = a'),
 Text(0.6046511627906976, 0.08333333333333333, 'gini = 0.0\nsamples = 648
\nvalue = [1057, 0]\nclass = a'),
 Text(0.813953488372093, 0.5833333333333334, 'NO_2 <= 39.5\ngini = 0.434\n
samples = 3589\nvalue = [3893, 1819]\nclass = a'),
 Text(0.7209302325581395, 0.4166666666666667, 'O_3 <= 23.5\ngini = 0.411\n
samples = 555\nvalue = [258, 634]\nclass = b'),
 Text(0.6744186046511628, 0.25, 'NO <= 15.5\ngini = 0.161\nsamples = 211\n
value = [30, 310]\nclass = b'),
 Text(0.6511627906976745, 0.08333333333333333, 'gini = 0.281\nsamples = 96
\nvalue = [27, 133]\nclass = b'),
 Text(0.6976744186046512, 0.08333333333333333, 'gini = 0.033\nsamples = 11
5\nvalue = [3, 177]\nclass = b'),
 Text(0.7674418604651163, 0.25, 'TOL <= 2.55\ngini = 0.485\nsamples = 344
\nvalue = [228, 324]\nclass = b'),
 Text(0.7441860465116279, 0.08333333333333333, 'gini = 0.32\nsamples = 209
```

# Conclusion

```python
In [35]: print("Linear Regression:",lr.score(x_test,y_test))
         print("Ridge Regression:",rr.score(x_test,y_test))
         print("Lasso Regression",la.score(x_test,y_test))
         print("ElasticNet Regression:",en.score(x_test,y_test))
         print("Logistic Regression:",logr.score(fs,target_vector))
         print("Random Forest:",grid_search.best_score_)
```

```
Linear Regression: 0.6259709949284429
Ridge Regression: 0.5877091721395884
Lasso Regression 0.2273508727383743
ElasticNet Regression: 0.3356452219334174
Logistic Regression: 0.9262454434993924
Random Forest: 0.935688248567957
```

# Random Forest Is Better!!!