

```
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_2013.csv")
```

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

Out[3]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL	station
0	2013-11-01 01:00:00	NaN	0.6	NaN	NaN	135.0	74.0	NaN	NaN	NaN	7.0	NaN	NaN	28079004
1	2013-11-01 01:00:00	1.5	0.5	1.3	NaN	71.0	83.0	2.0	23.0	16.0	12.0	NaN	8.3	28079008
2	2013-11-01 01:00:00	3.9	NaN	2.8	NaN	49.0	70.0	NaN	NaN	NaN	NaN	NaN	9.0	28079011
3	2013-11-01 01:00:00	NaN	0.5	NaN	NaN	82.0	87.0	3.0	NaN	NaN	NaN	NaN	NaN	28079016
4	2013-11-01 01:00:00	NaN	NaN	NaN	NaN	242.0	111.0	2.0	NaN	NaN	12.0	NaN	NaN	28079017

```
In [4]: df=df.fillna(1)
```

```
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'>
RangeIndex: 209880 entries, 0 to 209879
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        209880 non-null object
1   BEN         209880 non-null float64
2   CO          209880 non-null float64
3   EBE         209880 non-null float64
4   NMHC        209880 non-null float64
5   NO          209880 non-null float64
6   NO_2        209880 non-null float64
7   O_3         209880 non-null float64
8   PM10        209880 non-null float64
9   PM25        209880 non-null float64
10  SO_2        209880 non-null float64
11  TCH         209880 non-null float64
12  TOL         209880 non-null float64
13  station     209880 non-null int64
dtypes: float64(12), int64(1), object(1)
memory usage: 22.4+ MB
```

```
In [7]: data=df[['CO','station']]
data
```

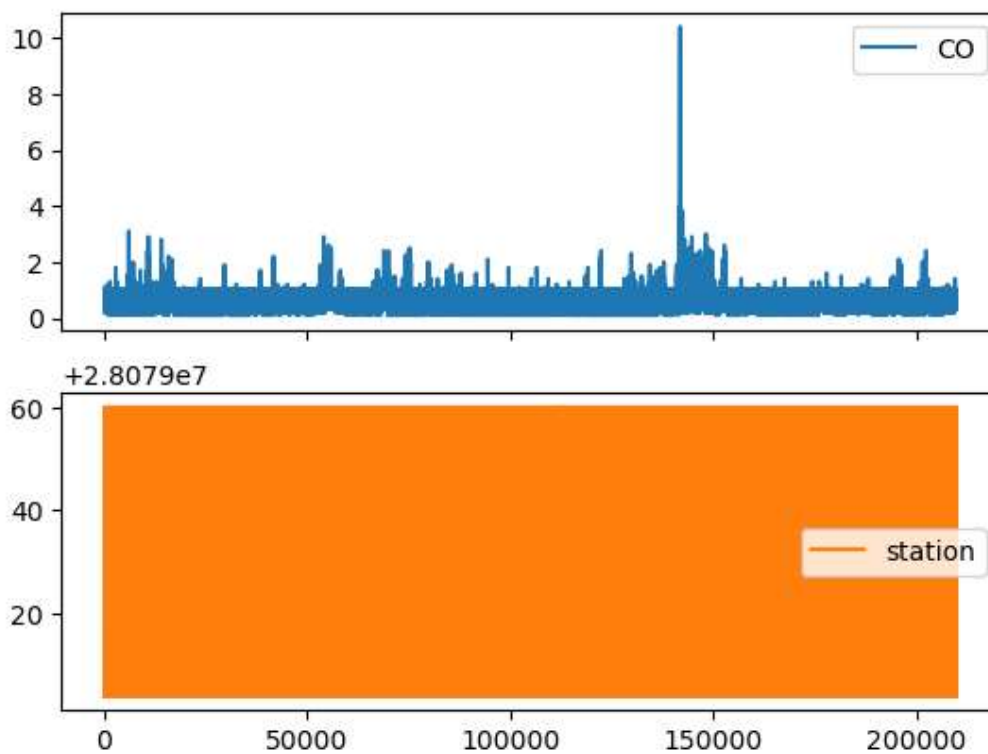
Out[7]:

	CO	station
0	0.6	28079004
1	0.5	28079008
2	1.0	28079011
3	0.5	28079016
4	1.0	28079017
...	...	...
209875	0.4	28079056
209876	0.4	28079057
209877	1.0	28079058
209878	1.0	28079059
209879	1.0	28079060

209880 rows × 2 columns

```
In [8]: data.plot.line(subplots=True)
```

Out[8]: array([<Axes: >, <Axes: >], dtype=object)



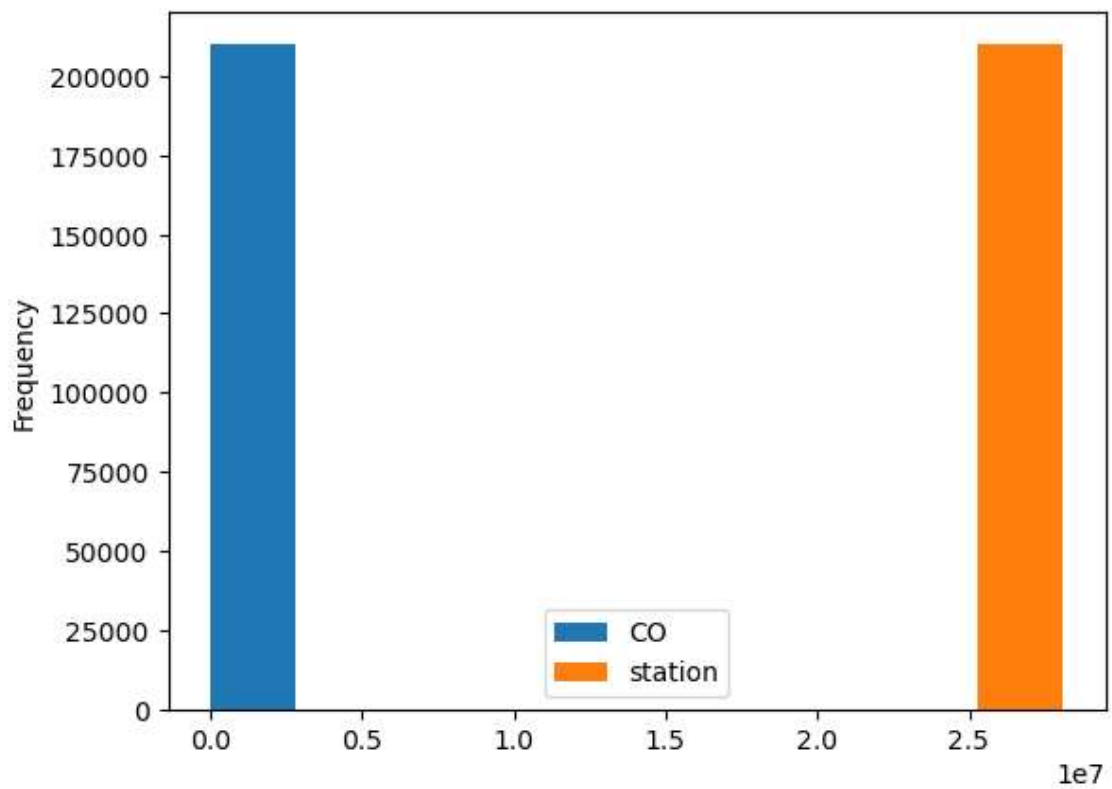
```
In [9]: b=data[0:50]  
b.plot.bar()
```

Out[9]: <Axes: >



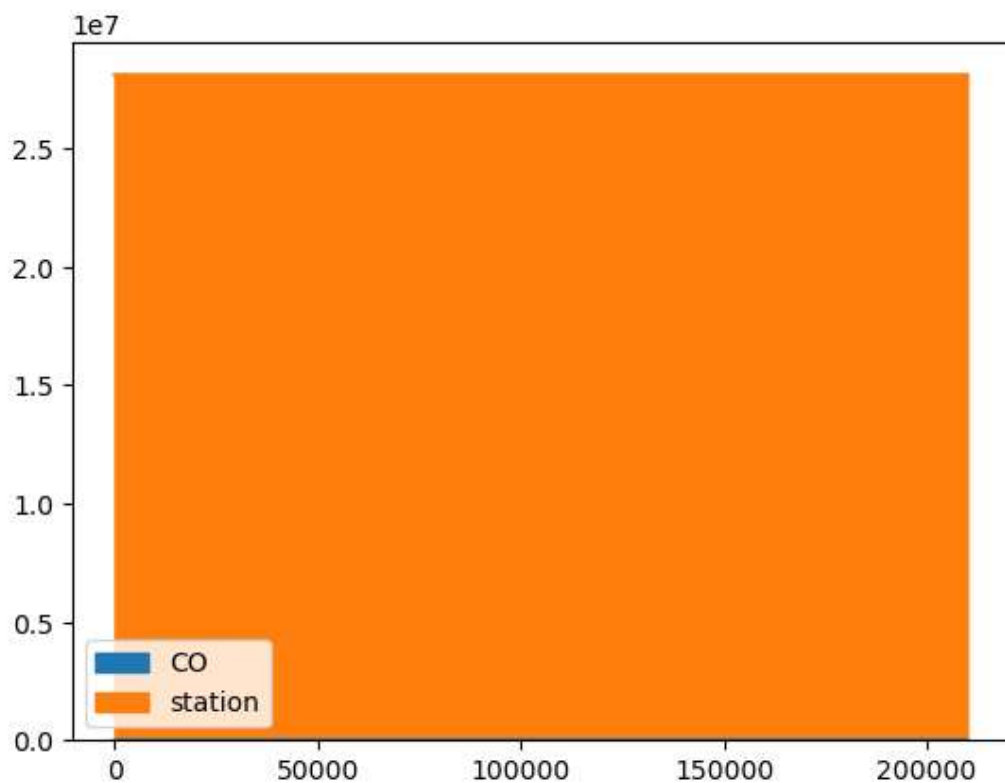
```
In [10]: data.plot.hist()
```

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



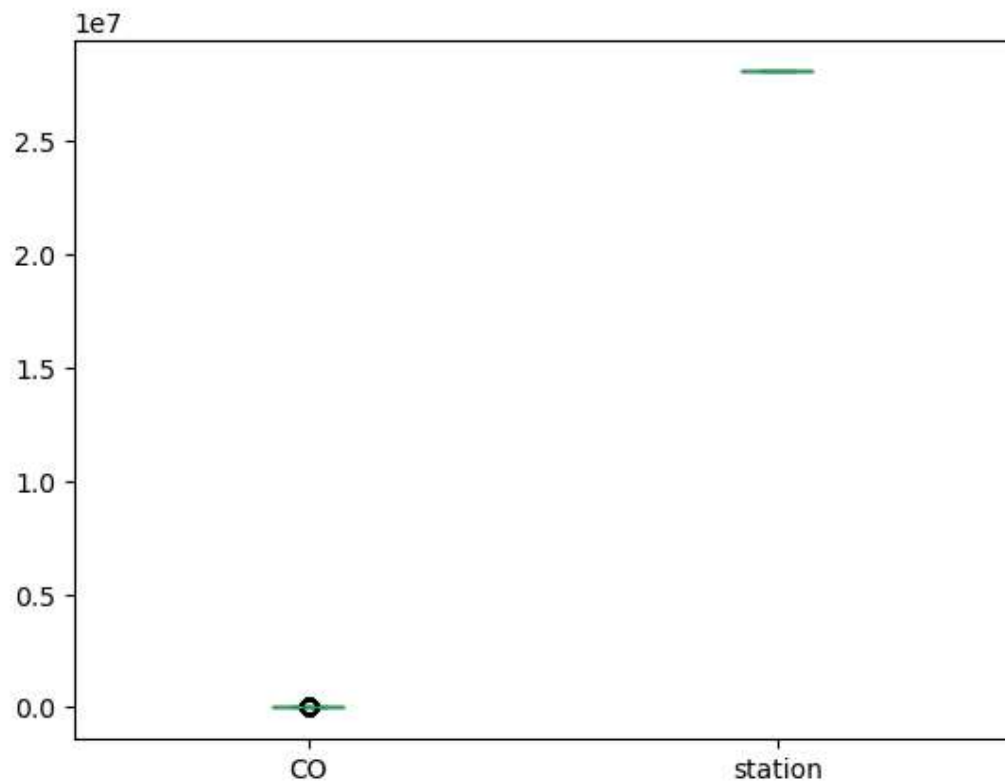
```
In [11]: data.plot.area()
```

```
Out[11]: <Axes: >
```



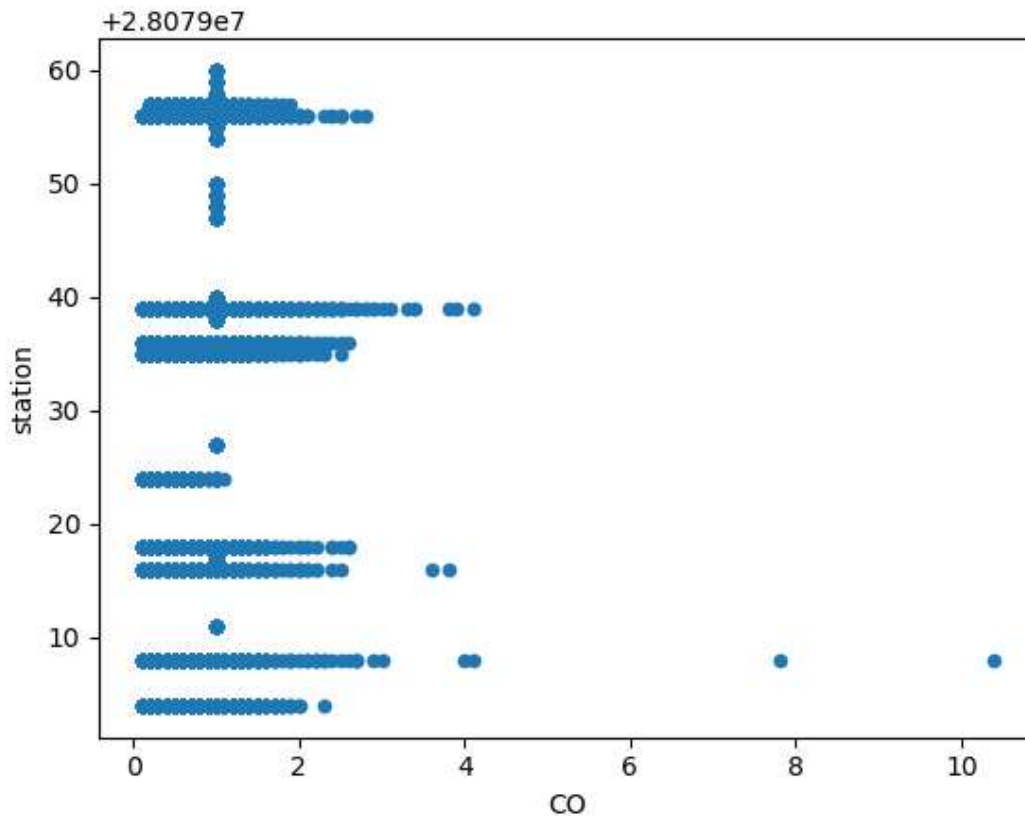
```
In [12]: data.plot.box()
```

```
Out[12]: <Axes: >
```



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

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



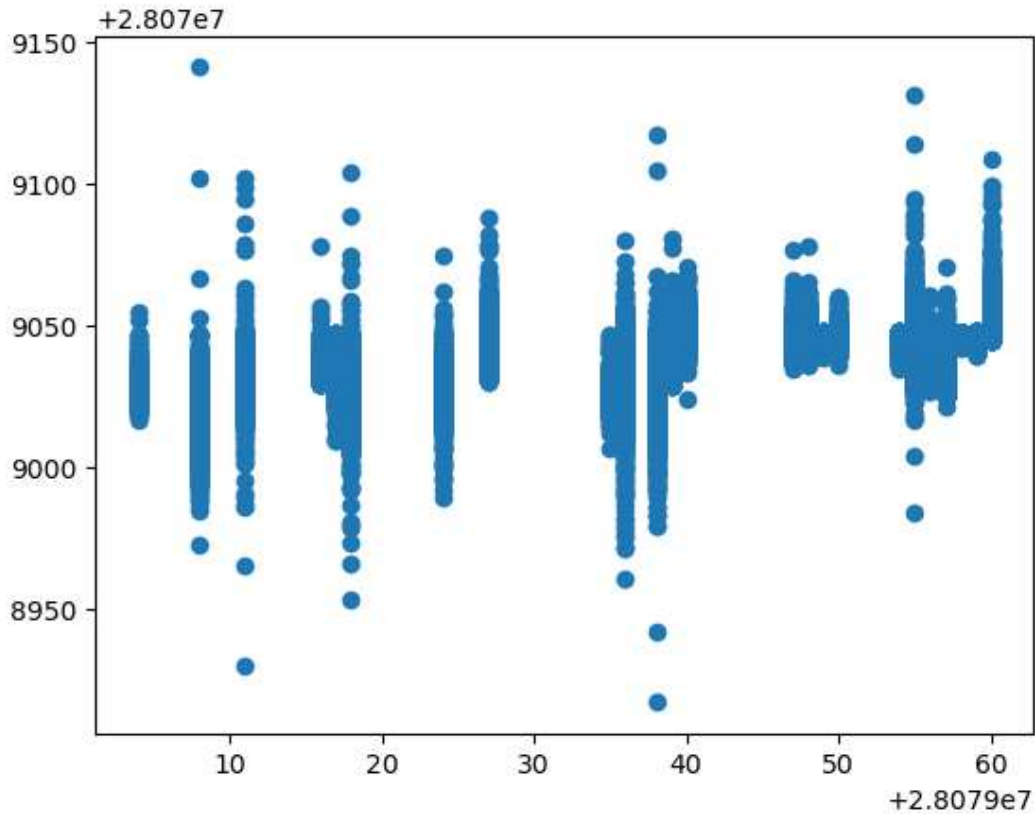
```
In [14]: x=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'NO', 'O_3',
               'PM10','PM25','SO_2', 'TCH', 'TOL']]
          y=df['station']
```

```
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]: 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 0x1f866faa050>



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

0.3041725824485335  
0.3085496347452269

## 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)
```

0.30417019167105386

0.3085467383822862

Out[18]:

▼	Lasso
Lasso(alpha=10)	

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

Out[19]: 0.04499070398821725

## 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([ 0.31266368, 2.6378251 , 0.47453002, 0. , -0.05521256,  
 0.03417485, -0.0204515 , 0.23453122, -0.35206307, -1.34674809,  
 -0. , -1.56991083])

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

Out[22]: 28079041.2491696

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

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

Out[24]: 0.16262799598578026

## 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)))
```

```
13.834458154386807
259.6173200826865
16.112644726508634
```

## 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']] [0:50]
target_vector=df[ 'station'] [0:50]
```

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

```
Out[29]: array([[1.32017385e-11, 1.77228839e-01, 3.04622063e-11, 6.75548895e-17,
5.48424019e-05, 1.81608859e-07, 1.09592662e-06, 1.98914252e-14,
2.81105441e-09, 1.18820886e-13, 8.22714249e-01, 4.75993568e-16,
1.90466882e-12, 6.60750110e-07, 7.34630911e-15, 1.26334810e-19,
2.00055463e-13, 6.55681339e-14, 1.28795666e-07, 3.79195253e-10,
4.64042021e-17, 3.25751097e-18, 1.25920588e-10, 2.87976550e-11]])
```

## Random Forest

```
In [30]: from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

```
Out[30]: RandomForestClassifier
RandomForestClassifier()
```

```
In [31]: parameters={'max_depth':[1,2,3,4,5],  
                  'min_samples_leaf':[5,10,15,20,25],  
                  'n_estimators':[10,20,30,40,50]  
                  }
```

```
In [*]: 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)
```

```
In [*]: 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','g','h','i','j','k','l','m','n','o','p','q','r','s','t','u','v','w','x','y','z'],f
```

## Conclusion

```
In [*]: 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_)
```

## logistic Regression Is Better!!!