

```
In [1]: import warnings
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split

warnings.filterwarnings("ignore")
warnings.filterwarnings("ignore", category=DeprecationWarning)

pd.set_option("display.max_rows", 999)
pd.set_option("display.max_columns", 590)
pd.set_option("display.width", 2000)
```

Data Exploration and key features

```
In [2]: # Invoking the prepared data
df = pd.read_csv('all_data.csv')

We see the features:

In [3]: df.sample()

Out[3]:
```

| | CO(GT)_winor | PT08.S1(CO2)_winor | PT08.S4(NO2)_winor | T_winor | RH_winor | AH_winor | day_winor | year_winor | quarter_winor | weekday_winor | NO2(GT)_winor | CO(GT)_winor | PT08.S1(CO2)_winor | CBH(GT)_winor | PT08.S2(NH4C)_winor | NO2(GT)_winor |
|------|--------------|--------------------|--------------------|----------|----------|----------|-----------|------------|---------------|---------------|---------------|--------------|--------------------|---------------|---------------------|---------------|
| 5991 | 0.572233 | 0.240188 | 0.52157 | 0.457463 | 0.501059 | 0.142594 | 0.758688 | 0.0 | 1.0 | 0.0 | 0.838242 | 2.4 | 925.0 | 8.6 | 923.0 | 442.0 |

We now consider the relevant features, i.e., such features that are normalized and transformed as we did in the `exploration` data notebook and the original data.

```
In [4]: key_features_original = [
    "CO(GT)",
    "PT08.S1(CO2)",
    "NO2(GT)",
    "PT08.S4(NO2)",
    "T",
    "RH",
    "AH",
    "day",
    "year",
    "quarter",
    "weekday",
]

key_features_winsor = [
    "CO(GT)_winsor",
    "PT08.S1(CO2)_winsor",
    "PT08.S4(NO2)_winsor",
    "NO2(GT)_winsor",
    "T_winsor",
    "RH_winsor",
    "AH_winsor",
    "day_winsor",
    "year_winsor",
    "quarter_winsor",
    "weekday_winsor",
]

In [5]: df_clusters = df[key_features_original+key_features_winsor]

In [6]: df_clusters.sample(3)

Out[6]:
```

| | CO(GT)_winor | PT08.S1(CO2)_winor | NO2(GT)_winor | PT08.S4(NO2)_winor | T | RH | AH | day | year | quarter | weekday | CO(GT)_winsor | PT08.S1(CO2)_winsor | PT08.S4(NO2)_winsor | NO2(GT)_winsor | T_winsor | RH_winsor | AH_winsor | day_winsor | year_winsor | quarter_winsor | weekday_winsor |
|------|--------------|--------------------|---------------|--------------------|------|------|--------|------|--------|---------|---------|---------------|---------------------|---------------------|----------------|----------|-----------|-----------|------------|-------------|----------------|----------------|
| 8016 | 1.9 | 1067.0 | 194.0 | 989.0 | 8.7 | 33.9 | 0.3832 | 7.0 | 2005.0 | 1.0 | 0.0 | 0.46965 | 0.441406 | 0.199005 | 0.039571 | 0.362982 | 0.423862 | 0.530810 | 0.505792 | | | |
| 4745 | 0.3 | 866.0 | 196.0 | 859.0 | 15.2 | 18.7 | 0.3214 | 21.0 | 2004.0 | 4.0 | 6.0 | 0.00000 | 0.147334 | 0.054884 | 0.600296 | 0.095458 | 0.035389 | 0.000000 | 0.874876 | | | |
| 6101 | 5.2 | 1387.0 | 104.0 | 2065.0 | 26.7 | 43.3 | 1.4881 | 22.0 | 2004.0 | 3.0 | 2.0 | 0.95966 | 0.811034 | 0.951642 | 0.039371 | 0.849216 | 0.658879 | 0.823688 | 0.891094 | | | |

K-means model

We first verify that our data satisfies normalization criteria:

```
In [7]: df_clusters[key_features_winsor].describe().T

Out[7]:
```

| | count | mean | std | min | 25% | 50% | 75% | max |
|---------------------|--------|----------|----------|-----|----------|----------|----------|-----|
| CO(GT)_winsor | 9357.0 | 0.458918 | 0.241618 | 0.0 | 0.291491 | 0.447019 | 0.609095 | 1.0 |
| PT08.S1(CO2)_winsor | 9357.0 | 0.545497 | 0.254899 | 0.0 | 0.264349 | 0.436113 | 0.631386 | 1.0 |
| PT08.S4(NO2)_winsor | 9357.0 | 0.565616 | 0.240293 | 0.0 | 0.431833 | 0.592653 | 0.723644 | 1.0 |
| NO2(GT)_winsor | 9357.0 | 0.589190 | 0.220316 | 0.0 | 0.448334 | 0.615927 | 0.727484 | 1.0 |
| T_winsor | 9357.0 | 0.629876 | 0.243586 | 0.0 | 0.491398 | 0.666862 | 0.824204 | 1.0 |
| RH_winsor | 9357.0 | 0.327853 | 0.255708 | 0.0 | 0.147486 | 0.676205 | 0.826590 | 1.0 |
| AH_winsor | 9357.0 | 0.519159 | 0.261592 | 0.0 | 0.345446 | 0.528677 | 0.715295 | 1.0 |
| day_winsor | 9357.0 | 0.704554 | 0.261762 | 0.0 | 0.548765 | 0.790067 | 0.906622 | 1.0 |
| year_winsor | 9357.0 | 0.240141 | 0.427182 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 1.0 |
| quarter_winsor | 9357.0 | 0.621679 | 0.382349 | 0.0 | 0.000000 | 0.442507 | 0.756471 | 1.0 |
| weekday_winsor | 9357.0 | 0.627278 | 0.326447 | 0.0 | 0.356207 | 0.712414 | 0.920782 | 1.0 |

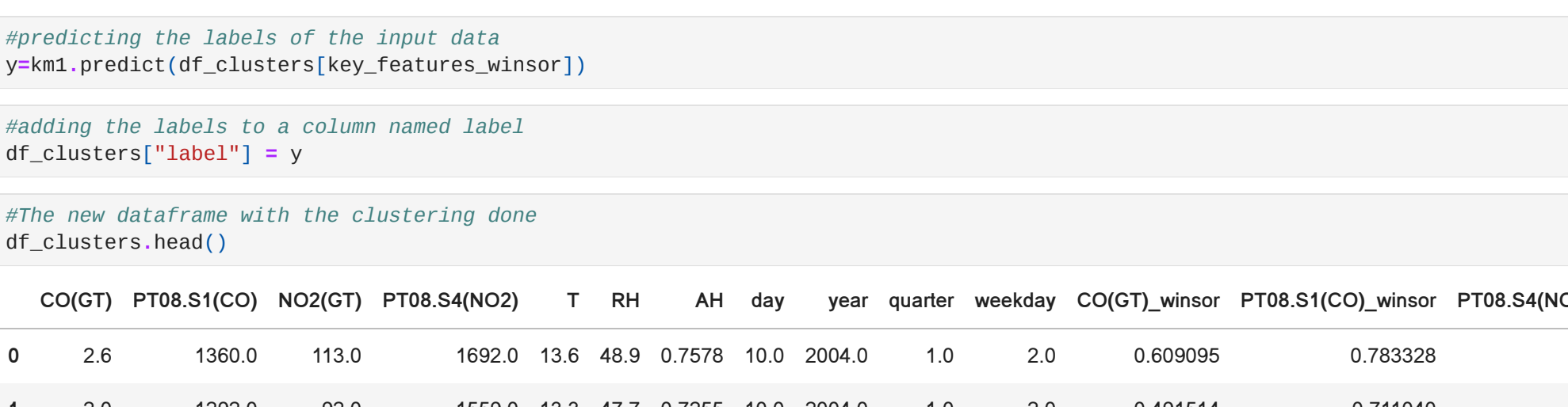
Since we are not dealing with other type of data types (besides of numeric), we then consider a K-means model:

```
In [8]: # Importing KMeans from sklearn
from sklearn.cluster import KMeans

In [9]: # Setting the number of clusters
KMEANS = 12

for i in range(1, KMEANS):
    km = KMeans(n_clusters=i, init='random', random_state=1492)
    km.fit(df_clusters[key_features_winsor])
    wcss.append(km.inertia_)

# The elbow curve
plt.figure(figsize=(12, 6))
plt.plot(range(1, KMEANS), wcss)
plt.plot(range(1, KMEANS), wcss, linewidth=2, color='red', marker='*')
plt.xlabel('% value')
plt.xticks(range(1, KMEANS, 1))
plt.ylabel('WCSS')
plt.show()
```



```
In [10]: # Making 5 clusters
km = KMeans(n_clusters=5, init='random', random_state=1492)
# Fitting the input data
km.fit(df_clusters[key_features_winsor])

Out[10]:
```

```
KMeans(init='random', n_clusters=4, random_state=1492)
```

```
In [11]: # Predicting the labels of the input data
y_km = km.predict(df_clusters[key_features_winsor])

In [12]: # Adding the labels to a column named label
df_clusters['label'] = y_km

In [13]: # The new dataframe with the clustering done
df_clusters.head()
```

| | CO(GT) | PT08.S1(CO2) | NO2(GT) | PT08.S4(NO2) | T | RH | AH | day | year | quarter | weekday | CO(GT)_winsor | PT08.S1(CO2)_winsor | PT08.S4(NO2)_winsor | NO2(GT)_winsor | T_winsor | RH_winsor | AH_winsor | day_winsor | year_winsor | quarter_winsor | weekday_winsor | label |
|---|--------|--------------|---------|--------------|------|------|--------|------|--------|---------|---------|---------------|---------------------|---------------------|----------------|----------|-----------|-----------|------------|-------------|----------------|----------------|-------|
| 0 | 2.6 | 1390.0 | 113.0 | 1692.0 | 13.6 | 48.9 | 0.3758 | 10.0 | 2004.0 | 1.0 | 2.0 | 0.60895 | 0.783328 | 0.747636 | 0.636718 | 0.546303 | 0.686762 | 0.354464 | 0.82136 | | | | 1 |
| 1 | 2.0 | 1390.0 | 92.0 | 1589.0 | 13.3 | 47.7 | 0.2585 | 10.0 | 2004.0 | 1.0 | 2.0 | 0.491514 | 0.711940 | 0.684323 | 0.521624 | 0.536803 | 0.655206 | 0.329419 | 0.65198 | | | | 1 |
| 2 | 2.2 | 1402.0 | 114.0 | 1555.0 | 11.9 | 54.0 | 0.1502 | 10.0 | 2004.0 | 1.0 | 2.0 | 0.533136 | 0.826194 | 0.681602 | 0.041056 | 0.487745 | 0.732847 | 0.348613 | 0.62188 | | | | 1 |
| 3 | 2.2 | 1370.0 | 122.0 | 1584.0 | 11.0 | 60.0 | 0.1967 | 10.0 | 2004.0 | 1.0 | 2.0 | 0.533136 | 0.799812 | 0.680497 | 0.079073 | 0.483938 | 0.803184 | 0.376487 | 0.62188 | | | | 1 |
| 4 | 1.6 | 1272.0 | 116.0 | 1450.0 | 11.2 | 59.6 | 0.1788 | 10.0 | 2004.0 | 1.0 | 2.0 | 0.399226 | 0.689053 | 0.617949 | 0.650062 | 0.461356 | 0.788715 | 0.378073 | 0.62188 | | | | 1 |

Analysis

```
In [14]: df_analysis = df_clusters[key_features_original+['label']]

In [15]: df_analysis.sample(3)

Out[15]:
```

| | CO(GT) | PT08.S1(CO2) | NO2(GT) | PT08.S4(NO2) | T | RH | AH | day | year | quarter | weekday | label |
|------|--------|--------------|---------|--------------|------|------|--------|------|--------|---------|---------|-------|
| 7393 | 0.7 | 1517.0 | 182.0 | 1802.0 | 14.1 | 67.8 | 0.5062 | 12.0 | 2005.0 | 1.0 | 2.0 | 0 |
| 4299 | 1.6 | 886.0 | 109.0 | 1341.0 | 27.1 | 26.9 | 0.9524 | 5.0 | 2004.0 | 3.0 | 6.0 | 2 |
| 5215 | 1.8 | 798.0 | 105.0 | 1123.0 | 15.0 | 56.4 | 0.9554 | 14.0 | 2004.0 | 4.0 | 3.0 | 2 |

```
In [16]: df_analysis.describe().T

Out[16]:
```

| | count | mean | std | min | 25% | 50% | 75% | max |
|--------------|--------|-------------|------------|-----------|-----------|-----------|-----------|----------|
| CO(GT) | 9357.0 | 2.089302 | 1.323034 | 0.1000 | 1.2000 | 1.8000 | 2.6000 | 11.900 |
| PT08.S1(CO2) | 9357.0 | 1096.382433 | 212.914565 | 647.0000 | 941.0000 | 1063.0000 | 1221.0000 | 2040.000 |
| NO2(GT) | 9357.0 | 112.373303 | 43.948519 | 2.0000 | 86.0000 | 109.0000 | 133.0000 | 340.000 |
| PT08.S4(NO2) | 9357.0 | 1456.528054 | 339.370078 | 551.0000 | 1242.0000 | 1463.0000 | 1662.0000 | 2775.000 |
| T | 9357.0 | 16.297574 | 8.705221 | -1.9000 | 12.0000 | 17.6000 | 24.1000 | 41.7000 |
| RH | 9357.0 | 49.248059 | 16.974949 | 9.3000 | 36.6000 | 49.6000 | 61.9000 | 88.700 |
| AH | 9357.0 | 1.024362 | 0.399578 | 0.1847 | 0.7461 | 0.9954 | 1.2962 | 2.201 |
| day | 9357.0 | 15.878884 | 8.808863 | 1.0000 | 8.0000 | 16.0000 | 23.0000 | 31.000 |
| year | 9357.0 | 2.004240141 | 0.427192 | 2004.0000 | 2004.0000 | 2004.0000 | 2004.0000 | 2005.000 |
| quarter | 9357.0 | 2.422571 | 1.543865 | 1.0000 | 1.0000 | 2.0000 | 3.0000 | 4.000 |
| weekday | 9357.0 | 3.000939 | 2.000323 | 0.0000 | 1.0000 | 3.0000 | 5.0000 | 6.000 |
| label | 9357.0 | 6.138666 | 1.122581 | 0.0000 | 1.0000 | 2.0000 | 3.0000 | 3.000 |

```
In [17]: df_analysis[df_analysis['label']==0].describe().T

Out[17]:
```

| | count | mean | std | min | 25% | 50% | 75% | max |
|--------------|--------|-------------|------------|-----------|-----------|-----------|-----------|----------|
| CO(GT) | 2247.0 | 2.054182 | 1.335288 | 0.1000 | 1.1000 | 1.8000 | 2.7000 | 8.700 |
| PT08.S1(CO2) | 2247.0 | 1107.112150 | 200.831137 | 715.0000 | 958.0000 | 1063.0000 | 1222.0000 | 1846.000 |
| NO2(GT) | 2247.0 | 141.070316 | 50.681205 | 17.0000 | 109.0000 | 138.0000 | 173.0000 | 340.000 |
| PT08.S4(NO2) | 2247.0 | 1163.969732 | 272.225393 | 551.0000 | 948.0000 | 1134.0000 | 1366.5000 | 2147.000 |
| T | 2247.0 | 10.495461 | 5.846498 | -1.9000 | 5.6000 | 9.8000 | 14.3000 | 30.000 |
| RH | 2247.0 | 62.262127 | 16.023949 | 9.9000 | 40.8000 | 50.9000 | 64.75000 | 86.000 |
| AH | 2247.0 | 0.621987 | 0.299947 | 0.1847 | 0.4463 | 0.6864 | 0.89735 | 1.193 |
| day | 2247.0 | 15.022897 | 8.916688 | 1.0000 | 7.0000 | 15.0000 | 23.0000 | 31.000 |
| year | 2247.0 | 2.005000000 | 0.000000 | 2005.0000 | 2005.0000 | 2005.0000 | 2005.0000 | 2005.000 |
| quarter | 2247.0 | 1.039878 | 0.102866 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 2.000 |
| weekday | 2247.0 | 3.003378 | 2.021764 | 0.0000 | 1.0000 | 3.0000 | 5.0000 | 6.000 |
| label | 2247.0 | 0.000000 | 0.000000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.000 |

```
In [18]: df_analysis[df_analysis['label']==1].describe().T

Out[18]:
```

| | count | mean | std | min | 25% | 50% | 75% | max |
|--------------|--------|-------------|------------|-----------|------------|-----------|------------|-----------|
| CO(GT) | 1466.0 | 1.607603 | 0.769205 | 0.100 | 1.100000 | 1.8000 | 1.90000 | 6.6000 |
| PT08.S1(CO2) | 1466.0 | 996.382356 | 147.832764 | 676.000 | 887.20000 | 968.0000 | 1088.0000 | 1520.000 |
| NO2(GT) | 1466.0 | 93.009550 | 32.392852 | 5.000 | 71.00000 | 101.0000 | 109.00000 | 239.0000 |
| PT08.S4(NO2) | 1466.0 | 1431.696499 | 242.877078 | 698.000 | 1286.00000 | 1439.0000 | 1591.00000 | 2381.000 |
| T | 1466.0 | 19.790109 | 7.928370 | 2.800 | 13.60000 | 19.1000 | 25.10000 | 41.7000 |
| RH | 1466.0 | 48.114734 | 17.284434 | 9.200 | 34.90000 | 48.0000 | 61.30000 | 82.3000 |
| AH | 1466.0 | 1.072660 | 0.367842 | 0.334 | 0.610075 | 1.0674 | 1.32575 | 2.0821 |
| day | 1466.0 | 17.047749 | 8.569119 | 1.000 | 10.00000 | 17.0000 | 25.00000 | 31.0000 |
| year | 1466.0 | 2.004000000 | 0.000000 | 2004.0000 | 2004.00000 | 2004.0000 | 2004.00000 | 2004.0000 |
| quarter | 1466.0 | 2.890588 | 0.917057 | 1.000 | 2.000000 | 3.0000 | 4.00000 | 4.0000 |
| weekday | 1466.0 | 0.390588 | 0.492864 | 0.000 | 0.000000 | 0.0000 | 1.000000 | 2.0000 |
| label | 1466.0 | 1.000000 | 0.000000 | 1.000 | 1.000000 | 1.0000 | 1.000000 | 1.0000 |

```
In [19]: df_analysis[df_analysis['label']==2].describe().T

Out[19]:
```

| | count | mean | std | min | 25% | 50% | 75% | max |
|--------------|--------|-------------|------------|-----------|------------|-----------|------------|-----------|
| CO(GT) | 3065.0 | 1.391648 | 0.567387 | 0.1000 | 0.9000 | 1.5000 | 1.8000 | 4.3000 |
| PT08.S1(CO2) | 3065.0 | 965.854160 | 110.334919 | 647.0000 | 883.0000 | 969.0000 | 1052.0000 | 1332.0000 |
| NO2(GT) | 3065.0 | 65.261223 | 28.371168 | 2.0000 | 43.0000 | 68.0000 | 109.0000 | 216.0000 |
| PT08.S4(NO2) | 3065.0 | 1387.815334 | 208.355866 | 682.0000 | 1279.0000 | 1438.0000 | 1529.0000 | 1932.0000 |
| T | 3065.0 | 19.790310 | 7.779521 | 1.2000 | 14.4000 | 19.1000 | 24.1000 | 42.2000 |
| RH | 3065.0 | 50.049060 | 18.677328 | 10.0000 | 38.0000 | 49.6000 | 62.3000 | 88.7000 |
| AH | 3065.0 | 1.029784 | 0.362103 | 0.1880 | 0.6831 | 0.9954 | 1.3974 | 2.1865 |
| day | 3065.0 | 16.146493 | 8.684626 | 1.0000 | 9.0000 | 16.0000 | 24.0000 | 31.0000 |
| year | 3065.0 | 2.004000000 | 0.000000 | 2004.0000 | 2004.00000 | 2004.0000 | 2004.00000 | 2004.0000 |
| quarter | 3065.0 | 2.900163 | 0.899929 | 1.0000 | 2.0000 | 3.0000 | 4.0000 | 4.0000 |
| weekday | 3065.0 | 4.319413 | 1.436231 | 2.0000 | 3.0000 | 5.0000 | 6.0000 | 6.0000 |
| label | 3065.0 | 2.000000 | 0.000000 | 2.0000 | 2.0000 | 2.0000 | 2.0000 | 2.0000 |

```
In [20]: df_analysis[df_analysis['label']==3].describe().T

Out[20]:
```

| | count | mean | std | min | 25% | 50% | 75% | max |
|--------------|--------|-------------|-------------|-----------|-----------|-----------|-----------|----------|
| CO(GT) | 2579.0 | 3.231601 | 1.436817 | 1.4000 | 2.1000 | 2.9000 | 4.0000 | 11.900 |
| PT08.S1(CO2) | 2579.0 | 1306.295463 | 179.938259 | 934.0000 | 1168.0000 | 1282.0000 | 1420.0000 | 2040.000 |
| NO2(GT) | 2579.0 | 130.610702 | 32.872330 | 27.0000 | 109.0000 | 122.0000 | 184.0000 | 278.000 |
| PT08.S4(NO2) | 2579.0 | 1807.203567 | 256.31212 | 1086.0000 | 1635.0000 | 1782.0000 | 1954.0000 | 2845.000 |
| T | 2579.0 | 22.427808 | 7.764322 | 1.4000 | 16.8000 | 21.8000 | 27.0000 | 44.600 |
| RH | 2579.0 | 46.315161 | 17.191819</ | | | | | |

Technical requirements

1. We would require MLOps frameworks to put in production such type of models.
2. It would be important to consider a machine learning architecture to integrate the different data sources with the different stages of the production phase.
3. Continuously monitoring of model performance and data distribution.
4. Real-time or near-to-real-time responses, depending of the final business goal. e.g. reduce the bad quality of air by locations, hours, etc.
5. Track the different models in order to manage the machine learning lifecycle.
6. A/B Testing to verify that different models in fact preserve suitable metrics in real-world scenarios.
7. Integration of different data sources: traffic, industrial outputs, and social media.