# Air pollution RATP AI mini-project 6

Team:
Lara Anna WAGNER (DiSc M2)
Yuliia NIKOLAENKO (DiSc M2)
Mihaela Elena GRIGORE (LeSc M2)

# Introduction to our dataset and problem

Quality of air monitoring:

- Underground:
    - From 3 subway stations (RATP)
    - Chatelet, Auber, Frank Roosevelt
    - Features:
        - NO (Nitrogen Monoxide), NO2 (Nitrogen Dioxide) - nitrogen oxides are produced from fuel combustion
        - CO2 - exhaled by commuters
        - PM10 (Particles suspended in air) - friction wheel vs rail, brake, construction material
        - Temperature and humidity

Purpose of gathering this data

- Health risks for commuters - something we must be aware of & estimate
- Ventilation system - project needs for new stations, evaluate satisfaction of needs in old stations, improve
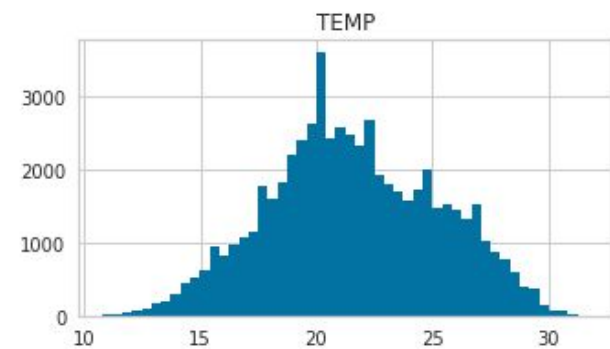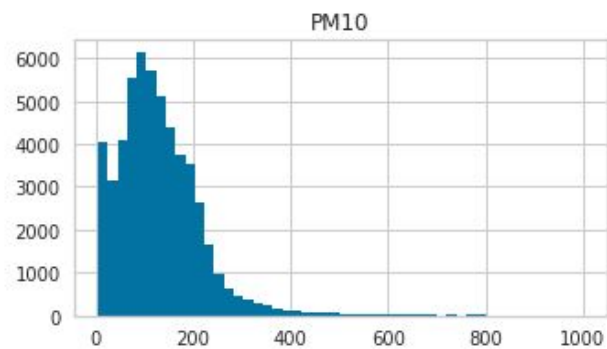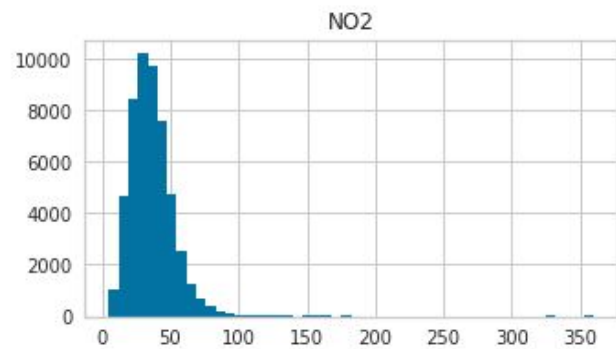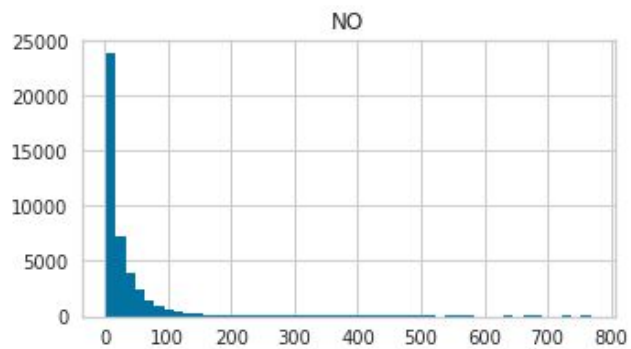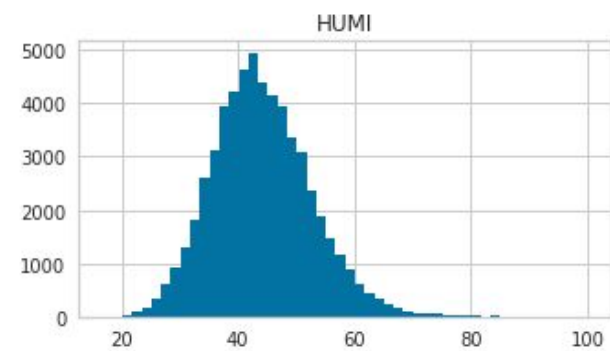- Are we replacing one source of pollution (surface) with another one (underground)
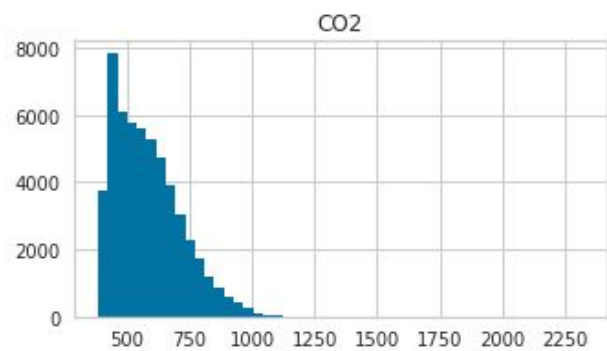
# Introduction to our dataset and problem
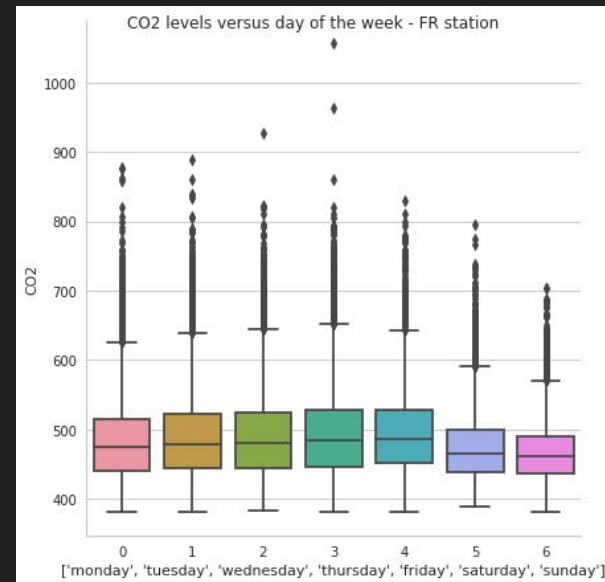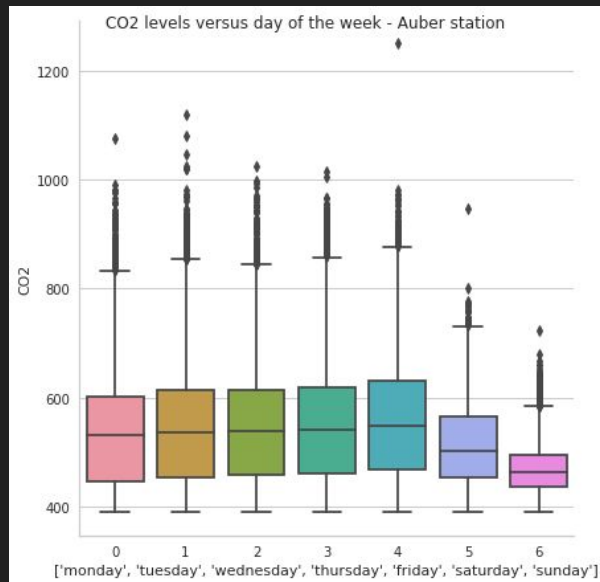
## EDA

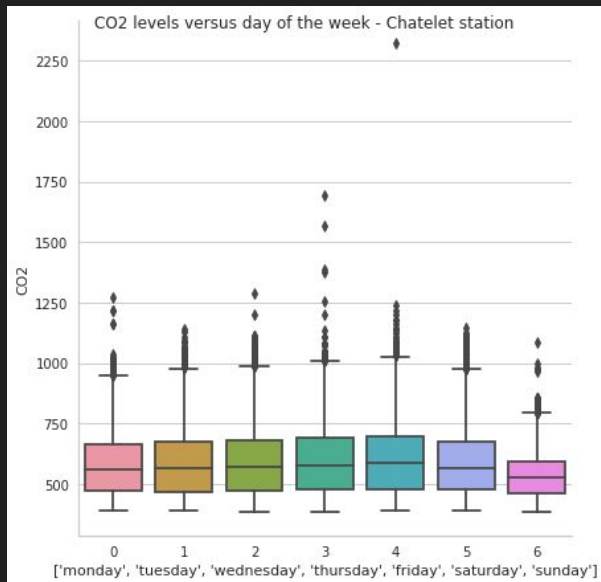- 7 years of data (from 2013 to Nov 2020)
- randomized
- 1 measurements per hour
- 24/7 measurements

## Data engineering

- Detecting and deleting missing values
- Converting the date into the right format: hour, weekday, month
- Creating new attributes weekend 0/1, summer 0/1

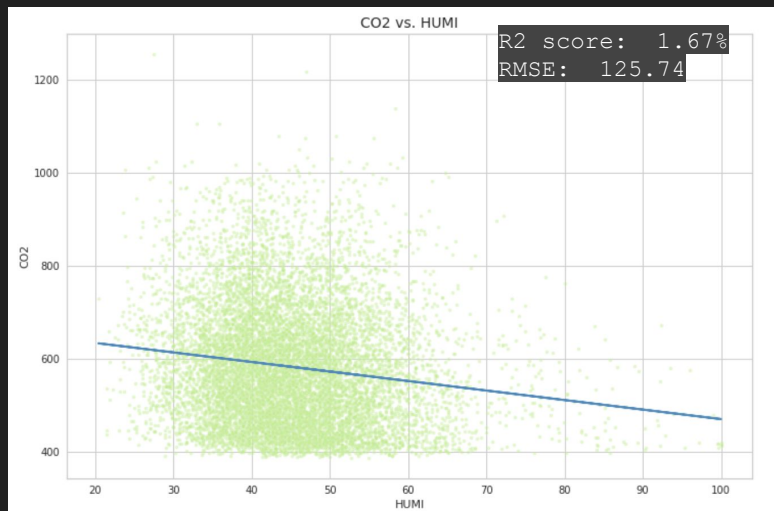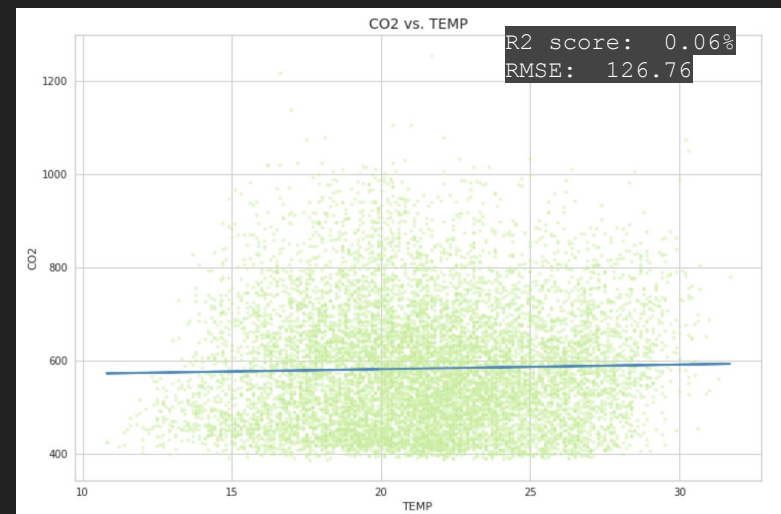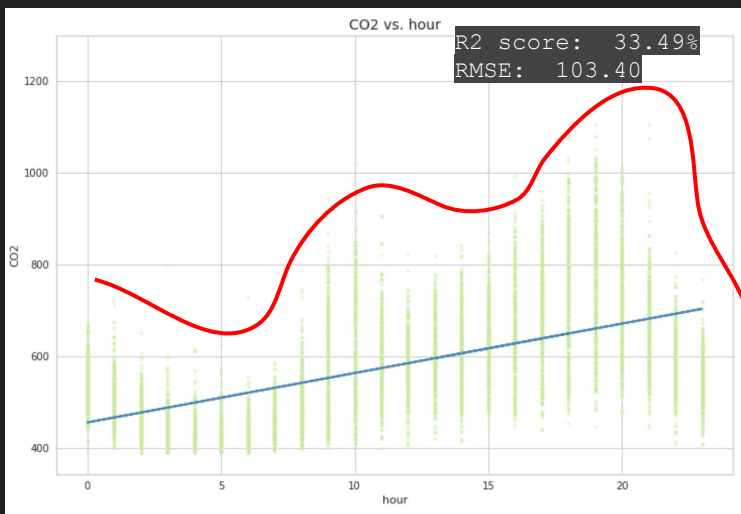# Weekday and CO2 level for each station

# Part 1

Predict **CO2** level in Chatelet from **time of day**, **temperature and humidity**

# Correlations



Indep:
TEMP, HUMI, Weekend, Hour

Dep: CO2

7

CO2 vs. hour
R2 score: 33.49%
RMSE: 103.40

CO2 vs. TEMP
R2 score: 0.06%
RMSE: 126.76

CO2 vs. HUMI
R2 score: 1.67%
RMSE: 125.74

| Feature | SD |
|---------|--------|
| CO2 | 128.52 |

# Evaluation CO2 in Chatelet

| | |
|---|---|
| Intercept | 582.79 |
| R2 Mean Score | 55% |
| RMSE | 80.45 |



Residuals for Ridge Model



Weight of selected features for predicting CO2 in Chatelete

Indep:
HUMI, Weekend, Hour
(ALL standardized)

Dep: CO2

# Part 2

Predict **NO2** level in a Chatelet from **past values**, **temperature** and **humidity + CO2**.

# Correlations



Indep:
TEMP, HUMI

Dep: NO2

| Feature | SD |
|---------|------|
| NO2 | 14.84 |

# Evaluation NO2 in Chatelet

Indep:
HUMI, TEMP

Dep: NO2

| Intercept | 36.17 |
|---|---|
| **R2 Mean Score** | 4.66% |
| **RMSE** | 14.39 |



Residuals for LinearRegression Model



Weight of selected features for predicting NO2 in Chatelete

| Feature | Coefficient |
|---|---|
| **TEMP** | -0.33 |
| **HUMI** | -3.92 |

13

# CO2 and NO2 correlation

# CO2 and NO2 correlation

# Comparison of the models predicting NO2 for Chatelet station

| Model | The Explained Variance | The Mean Absolute Error | The Median Absolute Error | RMSE | R2 score | R2 mean cross-valid |
|---|---|---|---|---|---|---|
| Without CO2 | 0.07 | 10.81 | 8.85 | 14.39 | 0.07 | 0.04 |
| With CO2 | 0.17 | 10.05 | 8.43 | 13.13 | 0.17 | 0.14 |

Better fit with the CO2 level as the the model shows higher accuracy and less errors probability.

# Part 3

Model Comparison

# Comparison of the models predicting CO2 and NO2 for all stations

| | Chatelet | | Auber | | Roosevelt | |
|---|---|---|---|---|---|---|
| Prediction of | CO2 | NO2 | CO2 | NO2 | CO2 | NO2 |
| The Explained Variance | 0.24 | 0.32 | 0.31 | 0.38 | 0.43 | 0.38 |
| The Mean Absolute Error | 85.72 | 9.16 | 114.31 | 15.75 | 33.2 | 13.2 |
| The Median Absolute Error | 70.22 | 7.65 | 108.79 | 12.75 | 70.22 | 7.65 |
| RMSE | 112.16 | 11.92 | 137.08 | 21.23 | 44.59 | 17.23 |
| R2 score | 0.24 | 0.32 | -0.62 ? | -0.1 | 0.43 | 0.38 |
| R2 mean cross-valid | 0.2 | 0.3 | 0.3 | 0.37 | 0.4 | 0.35 |

# Part 4

Based on data from two stations, can we predict air quality in the third one ?

# Revisit correlations - what features to consider

# The new dataframe

[feat_1_station_1, …, feat_n_station_1, feat_1_station_2, …., feat_n_station_2, dependent_variable_1_station_3]

We inner join df_chatelet, df_auber and df_roosevelt.CO2 on the DATE/HEURE colum =>

- NO
- NO2
- PM10
- TEMP      one set for each station (Chatelet and Auber)
- HUMI
- Hour      0:23
- Weekday      0:6
- Month      1:12
- Weekend      0/1
- Summer      0/1

# Wrangling and feature engineering

- We remove all rows with at least one missing value.
- From ~68.000 entries, we are left with ~20.000

# Outliers

- Remember we saw outliers in previous slides
- Let's see what remained after we removed ⅔ of our data (na removal)



Box plots for CO2, NO, NO2, PM10 - Chatelet dataset

# MLR model coeff



Independent
NO, NO2, PM10,
TEMP, HUMI

Hour, Weekday, Month,
Weekend, Summer

ALL standardized

Chatelet + Auber

Dep: CO2 F.R. station

https://stats.stackexchange.com/questions/463690/multiple-regression-with-mixed-continuous-categorical-variables-dummy-coding-s

# MLR model - evaluation



The Explained Variance: 0.56

The Mean Absolute Error: 30.33

The Median Absolute Error: 23.46

Mean squared error: 1652.41

Root mean squared error: 40.65

R2 Score: 0.56

Intercept: 490.64

# The Effect of regularization

# New dataframe: one hot encoding

We inner join df_chatelet, df_auber and df_roosevelt.CO2 on the DATE/HEURE colum =>

- NO
- NO2
- PM10
- TEMP          one set for each station (Chatelet and Auber)
- HUMI
- Hour          0:23 ---|
- Weekday       0:6   ---|---> categorical variables or not ? → one hot encoding
- Month         1:12 ---|
- Weekend       0/1
- Summer        0/1

# New dataframe: one hot -> exploding coefficients



Weight of each feature for predicting CO2 in F.R. station

NO
NO2
PM10
TEMP
HUMI
Hour        0:23 ---|
Weekday     0:6   ---|-one-hot
Month       1:12 ---|
Weekend     0/1
Summer      0/1

ALL standardized

Even the one hot encoded variables

Conclusion: must prevent large coefficients

28

# Ridge regularization prevents large coefficients



Weight of each feature for predicting $CO_2$ in F.R. station

Observations:
- We don't find the same relationship between increasing hour and coeff size / sign
- Difficult to interpret
- We try Lasso next !

```
The Explained Variance: 0.57
The Mean Absolute Error: 29.81
The Median Absolute Error: 23.62
Mean squared error: 1579.47
Root mean squared error: 39.74
```

# Lasso regularization -> to reduce # of features



Weight of each feature for predicting CO2 in F.R. station

alpha = 0.8 (where alpha = 0 is equivalent to ordinary least square)

Observations:
- CO2 in Chatelet and Auber most important features when using Lasso regularization
- We sacrificed performance ?

```
The Explained Variance: 0.56
The Mean Absolute Error: 30.24
The Median Absolute Error: 23.96
Mean squared error: 1635.66
Root mean squared error: 40.44
R2 Score: 0.56
Intercept: 490.622075
```
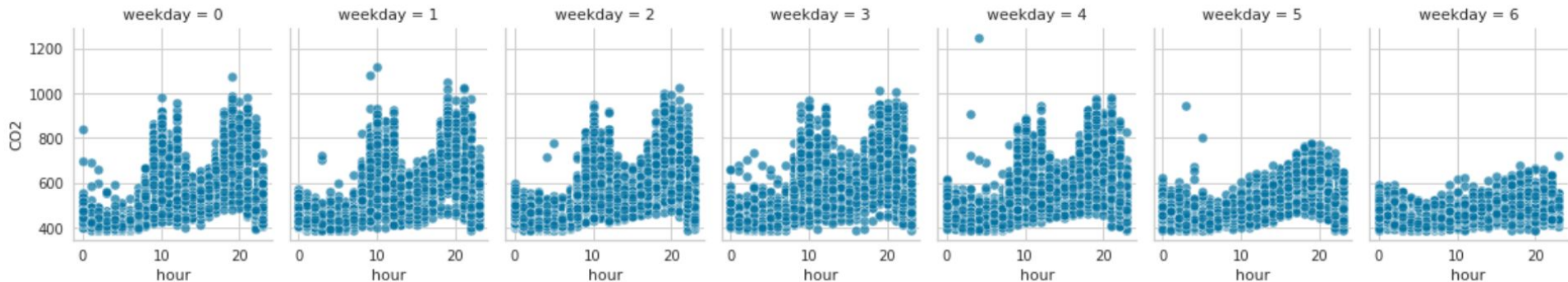
# Non - linear relationship - CO2 versus hour per day

- We see here why hour of day (0:23) would not be a good candidate to predict $CO_2$ on through a linear model



CO2 levels versus hour for each day of the week - Auber station

# Models comparison

# Model comparison: same features,different approach

Independent: NO, NO2, PM10, TEMP, HUMI, hour, weekday, month, weekend, summer in Chatelet & Auber
Dependent: $CO_2$ in F. Roosevelt
Model 1: all features numeric, scaling, no regularization
Model 2 : all features, one hot encoding + scaling + Ridge regularization
Model 3: all features, one hot encoding + scaling + Lasso regularization
Question: differences in outcome ?

|          | Model 1 | Model 2 | Model 3 |
|----------|---------|---------|---------|
| R2 score | .56     | .57     | .56     |
| RMSE     | 40.65   | 39.74   | 40.44   |

# Model comparison: 5-cv

Independent: NO, NO2, PM10, TEMP, HUMI, hour, weekday, month, weekend, summer in Chatelet & Auber
Dependent: CO2 in F. Roosevelt
Model 1: all features, one hot encoding + scaling no regularization
Model 2 : all features, one hot encoding + scaling + Ridge regularization
Model 3: all features, one hot encoding + scaling + Lasso regularization
Question: differences in outcome ?

```
R2 mean score simple lr: 0.580
R2 mean score ridge: 0.580
R2 mean score lasso: 0.565
```

# Next

RFE to choose features

Data engineering:

- Hour and month as sine (keep cyclical nature)
- Or encode hour as rush hour / not rush
- Imputation for missing values



CO2 levels versus hour for each day of the week - Auber station