# Acute Respiratory Infections Forecast using Machine Learning
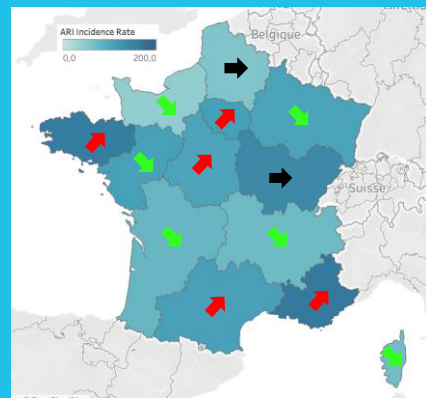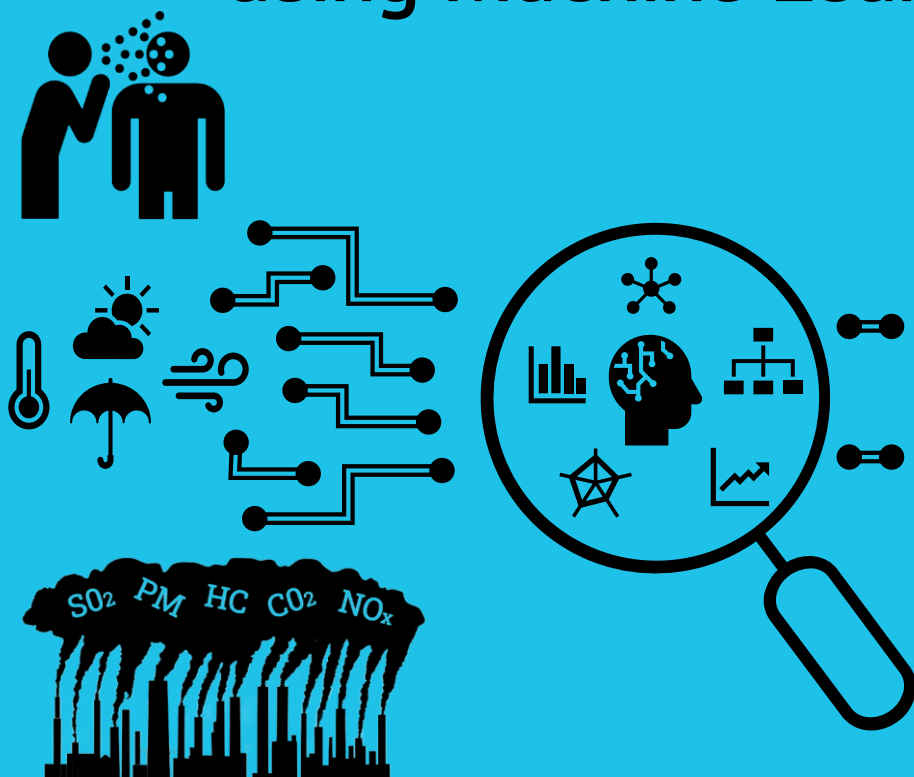
Ironhack Data Analyst certification project
Romain Courtois

# 1. Case Study

**A**cute **R**espiratory **I**nfections :

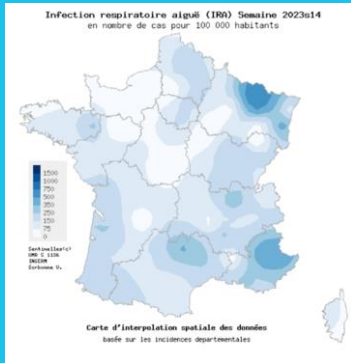Caused by various respiratory viruses including SARS-CoV-2

Study focused on factors of the disease: pollution, climate and seasons.

Using **ETL**, **EDA** and **ML**

→ build a model that can predict future **incidence rate**

# 2. Data Sources



## A. French epidemiological surveillance:
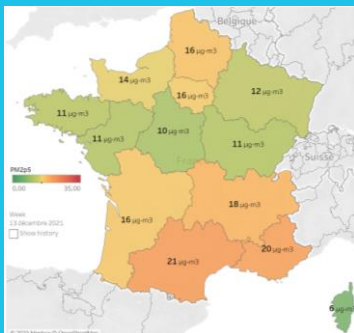
Weekly ARI incidence rate by regions



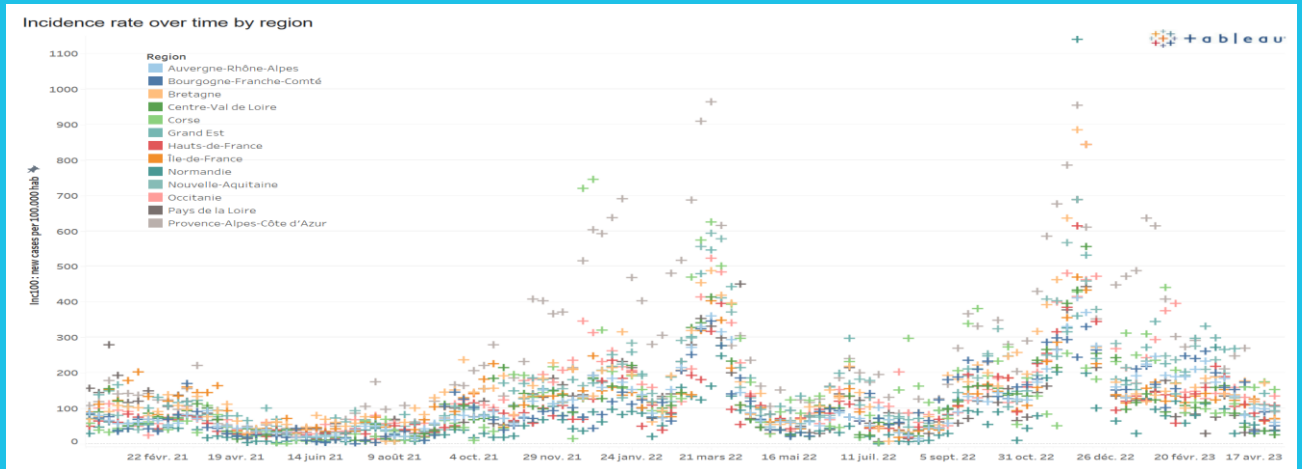## B. Historical meteorological observation France:

Daily temp., pressure, humidity…





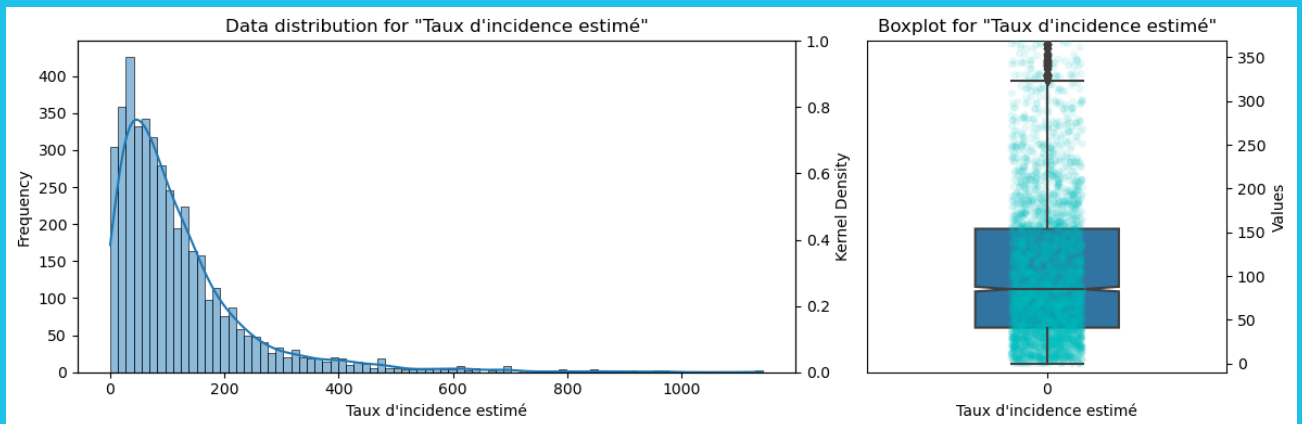## C. Concentrations of air pollutants

O3, NOx, SO2, PM10, PM2.5, CO, CH6H6

# 3. Extract Transform Load

**Extract**

API → aggregated data

.CSV → raw data

**Transform**

**Python :**

aggregate data by weeks

normalize categories/regions

**Load**

**MySQL** relational database
- Staging area
- DataMart

# 4. Exploratory Data Analysis
## Incidence rate over time



Incidence rate over time by region

## Data distribution

# 5. Iterative Machine Learning

**A. Data manipulations**

**B. Features selection**

**E. use feedback as new parameters**

scikit learn

**C. Train regressions models**

**D. Evaluate best models**
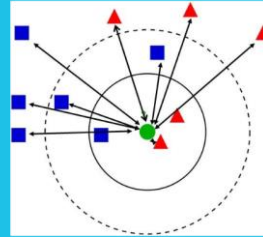
**Use ML to take decisions**

# 6. Regression models used

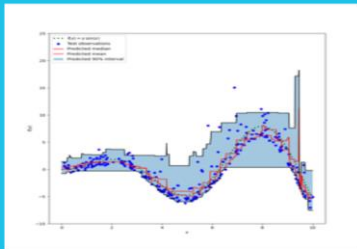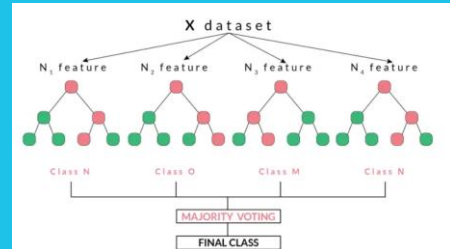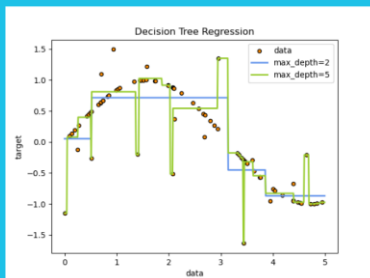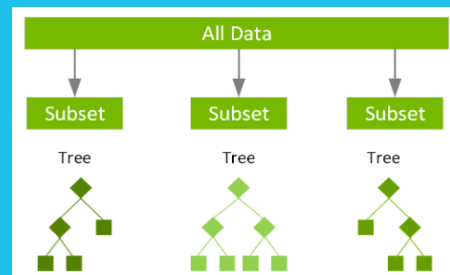- Linear Regression



- K Neighbors



- Gradient Boosting



- Random Forest
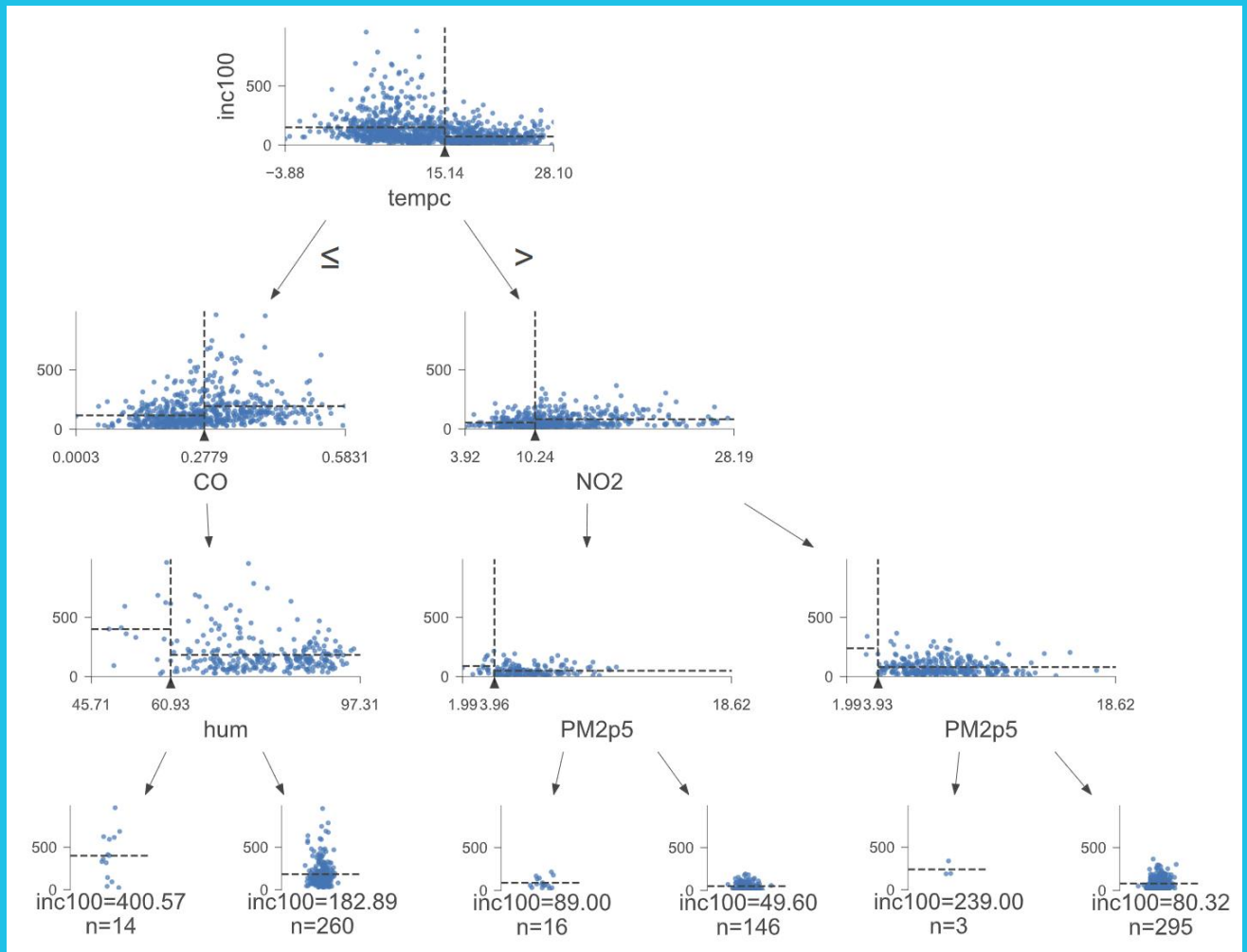


- Decision Tree



- XGBoost

# 7. Focus on XGBoost tree regression

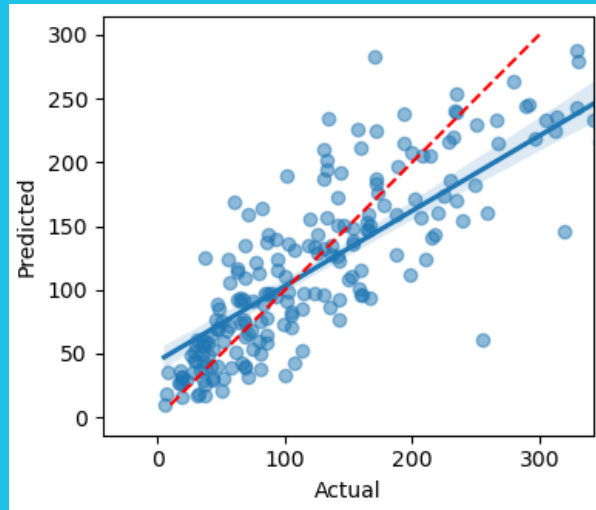

CO and humidity is used when temp < 15°c
NO2 and PM2.5 when >15°c

# 8. Conclusion

Best prediction score : Gradient Boosting
37% (mape)



With more data, it improves over time.

I will publish the prediction dashboard when score is < 10%.

Follow me on :

**github.com/romaincrt/ARIForecast**