# Air Quality 🌨️

1. **Introduction**
2. **Collect Data**
3. **Clean Data**
4. **Processing Data**
5. **Analysis**
6. **Conclusions**

# Introduction

### Introduction

- Big Data : Powerful method to analyze for insights from the sporadic data.
- We aim to increase the knowledge on data analysis from this project

### Project Subject

- Air quality of the lleida city between 2019 and 2020
- How the pandemic situation affect to air quality of Lleida city?
- Based on abundance in Microgram per Cubic Meter of Air (ug/m3) of $O_3$, $CO$, $NO_2$, $SO_2$

# Introduction

### Environment

- Python3
- Anaconda : python virtual environment
- Docker : spark context
- Jupyter Notebook , Python pyspark module
- Cookiecutter : project structure

# Data Collecting

**Source**

- Open AQ API : non-profit organization empowering communities around the globe

**Data collecting**

- Building a python script "get_data.py" via HTTP request **OpenAQ API, JSON format**
- **Parameters.json :** units & types - **$O_3$, CO, $NO_2$, $SO_2$**
- **Locations.json :** sensor information in Lleida city
- **Measurement.json :** the obtained data from each sensor

**Issues**

- Choosing Source : different models, categorization and formats by each data source
- Comparing the data of the city : limited computational power

# Data Cleaning

- Jupyter notebook
- Pandas Profiling => Reports
- Remove  useless data
- Remove  negative values from the sensors => lack of information
- Re-design initial model

# Data Cleaning

## locations schema

```
{
    "city": "string",
    "country": "string",
    "measurements": "int",
    "name": "string",
    "parameters": [{
        "average": "float",
        "count": "int",
        "displayName": "string",
        "firstUpdated": "datetime",
        "id": "int",
        "lastUpdated": "datetime",
        "lastValue": "int",
        "parameter": "string",
        "parameterId": "int",
        "unit": "string"
    }]
}
```

## parameters schema

```
{
    "displayName": "string",
    "name": "string",
    "preferredUnit": "string"
}
```

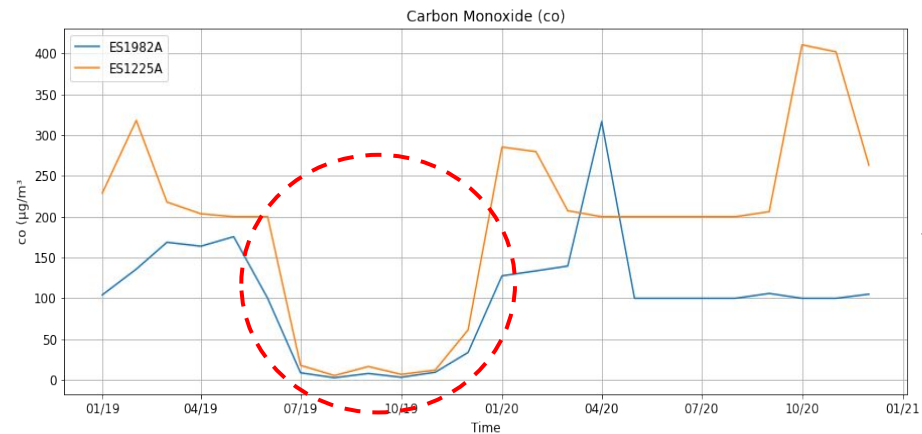## measurements schema

```
{
    "location": "string",
    "city": "string",
    "date": {
        "local": "datetime",
        "utc": "datetime"
    },
    "parameter": "string",
    "value": "float",
    "unit": "string"
}
```
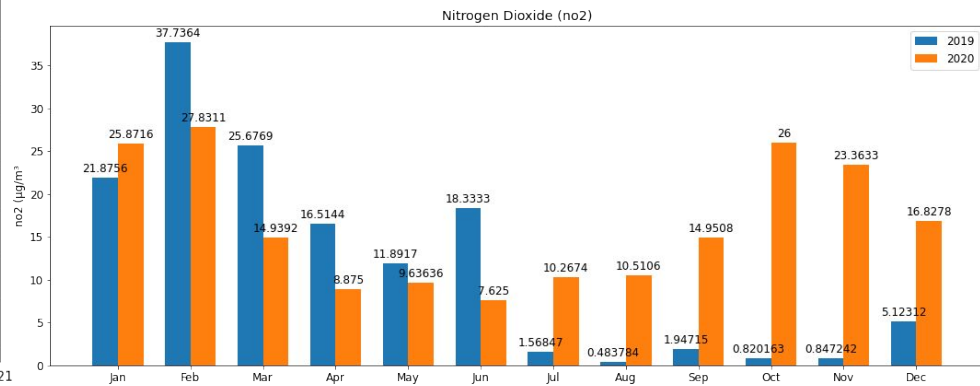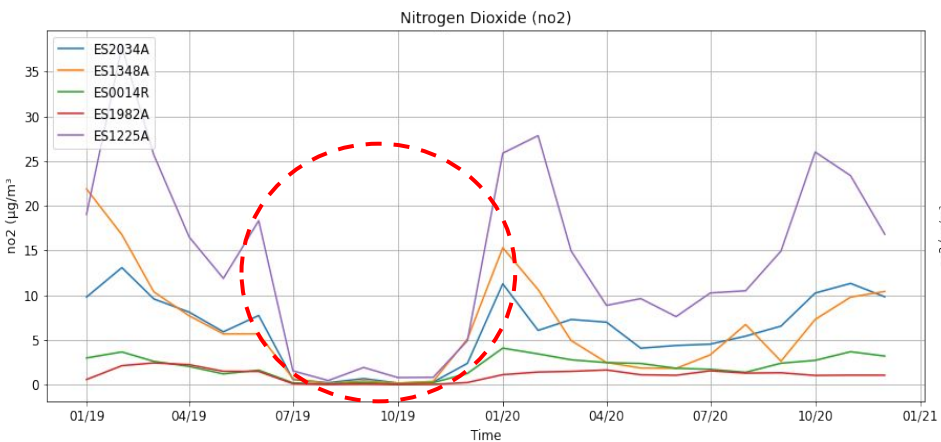
# Data Processing

- Docker Image: **jupyter/all-spark-notebook -> pyspark + notebook + matplotlib**
- Group measurements by: **Sensor**
- Aggregate: Date by **month**
- **JOIN ALL**
- Parameters: **NO2, CO, SO2, O3**
- **! Issue ! - No available data** during the last six months of **2019**
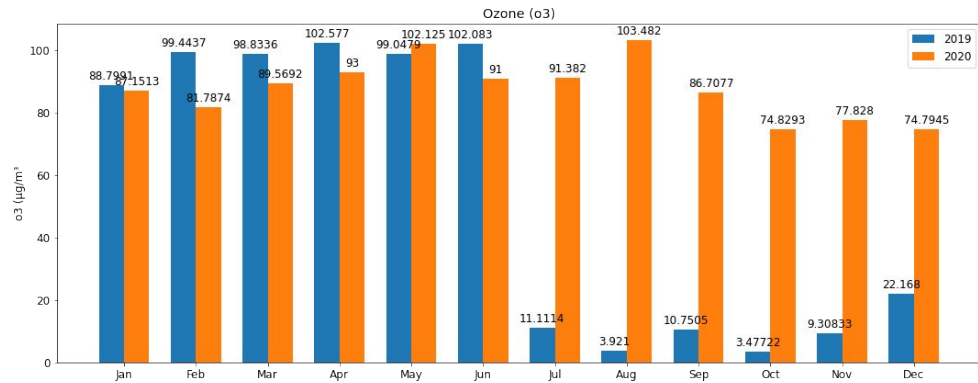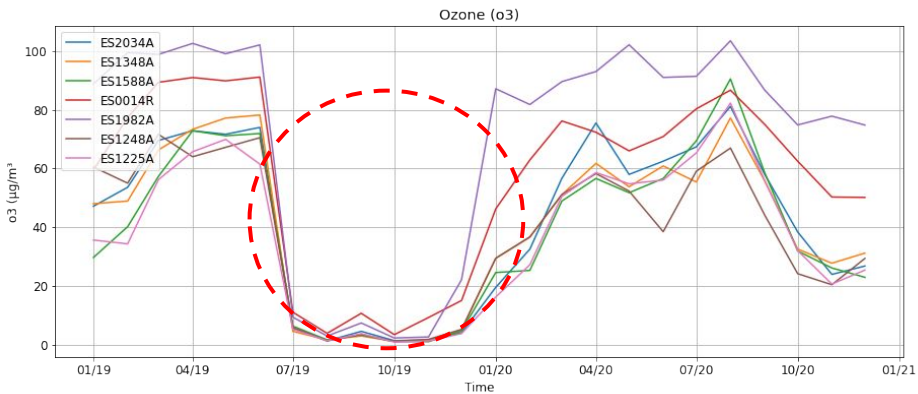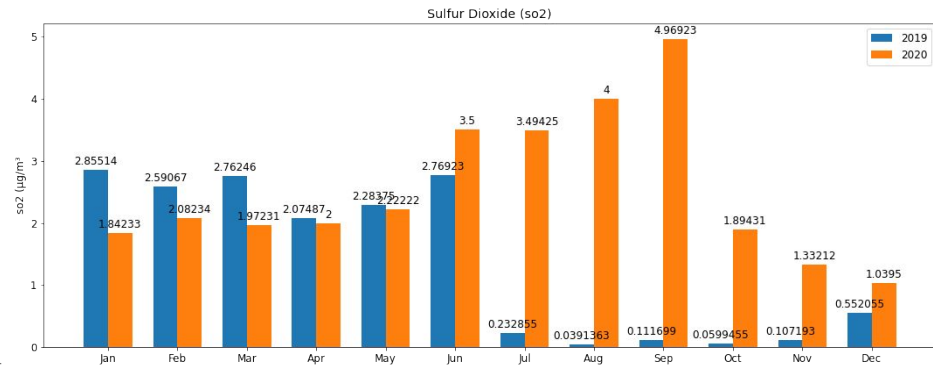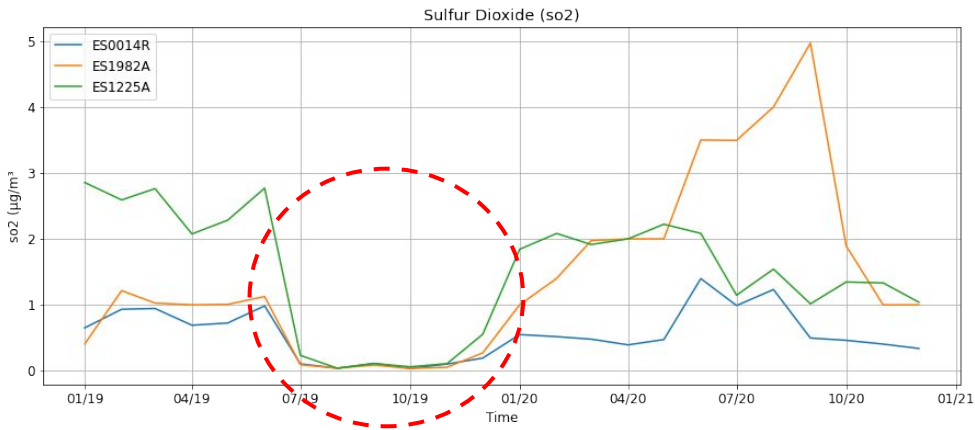- Air Sensors:  **High volume of data + different precision + location**

# Analysis (CO)



Carbon Monoxide (co)



Carbon Monoxide (co)

# Analysis (NO2)

# Analysis (O3)

# Analysis (SO2)



Sulfur Dioxide (so2)

# Conclusions