

DATA ANALYSIS FOR US AIR POLLUTION

XIANG LIU
XUE XU

AGENDA

- Project Overview
- Inspirations
- Data Pre-processing
- Data Analysis
- Conclusion
- Reference

PROJECT OVERVIEW

This dataset includes four major pollutants (Nitrogen Dioxide, Sulphur Dioxide, Carbon Monoxide and Ozone) for every day from 2000 - 2016 and place them neatly in a CSV file

Pollutants	Cause
NO ₂	are expelled from high temperature combustion, and are also produced during thunderstorms by electric discharge
SO ₂	produced by volcanoes and in various industrial processes
O ₃	brought about by human activities (largely the combustion of fossil fuel)
CO	vehicular exhaust contributes to the majority of carbon monoxide let into our atmosphere

TOXITY

Pollutants	Toxity
NO ₂	Breathing difficulties, Throat spasms, Headache, Fatigue, Nausea
SO ₂	Temporary respiratory problems, chronic bronchitis, emphysema, decreased fertility, coughing
O ₃	Irritant effect on the respiratory tract; stimulation of atherosclerosis
CO	Confusion, vision and balance problems, Loss of consciousness, Nausea and vomiting, Headaches

DATASET

There is a total of 28 fields. The four pollutants (NO₂, O₃, SO₂ and CO) each has 5 specific columns. Observations totaled to over 1.4 million

Column Name	Description
NO ₂ Units	The units measured for NO ₂
NO ₂ Mean	The arithmetic mean of concentration of NO ₂ within a given day
NO ₂ AQI	The calculated air quality index of NO ₂ within a given day
NO ₂ 1st Max Value	The maximum value obtained for NO ₂ concentration in a given day
NO ₂ 1st Max Hour	The hour when the maximum NO ₂ concentration was recorded in a given day

AQI --- AIR QUALITY INDEX

Air Quality Index Levels of Health Concern	Numerical Value	Meaning
Good	0 to 50	Air quality is considered satisfactory, and air pollution poses little or no risk.
Moderate	51 to 100	Air quality is acceptable; however, for some pollutants there may be a moderate health concern for a very small number of people who are unusually sensitive to air pollution.
Unhealthy for Sensitive Groups	101 to 150	Members of sensitive groups may experience health effects. The general public is not likely to be affected.
Unhealthy	151 to 200	Everyone may begin to experience health effects; members of sensitive groups may experience more serious health effects.
Very Unhealthy	201 to 300	Health warnings of emergency conditions. The entire population is more likely to be affected.
Hazardous	301 to 500	Health alert: everyone may experience more serious health effects.

DATASET

Column Name	Description
State Code	The code allocated by US EPA to each state
County Code	The code of counties in a specific state allocated by US EPA
Site Num	The site number in a specific county allocated by US EPA
Address	Address of the monitoring site
State	New Jersey...
County	County of monitoring site
City	City of the monitoring site
Date	LocalDate of monitoring

INSPIRATIONS

1. How is the concentration of each pollutant in US states?
2. Which time in a day has the maximum density of air pollutants?
3. From 2000-2016, the concentration of each pollutant increase or decrease?
4. Can we predict the trend of each pollutant in the future?

DATA PREPROCESSING

Step 1 check if the data has missing value

```
In [3]: le=data['CO AQI']
le[le.isnull()]
```

```
Out[3]: 0      NaN
         2      NaN
         4      NaN
         6      NaN
         8      NaN
        10     NaN
        12     NaN
```

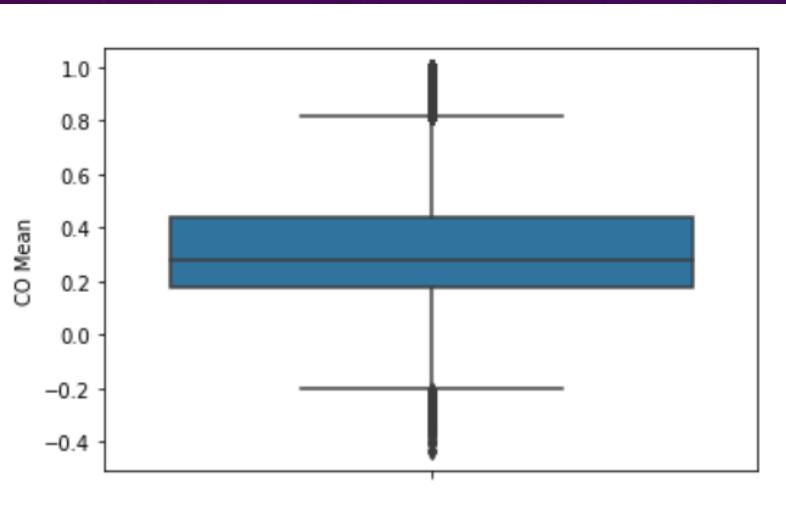
```
data = pd.read_csv('/Users/xiangliu/Desktop/CSC560 Data/pollution_us_2000_2016.csv')
data.shape
(1746661, 29)
```

Step 2 Use dropna to drop missing value

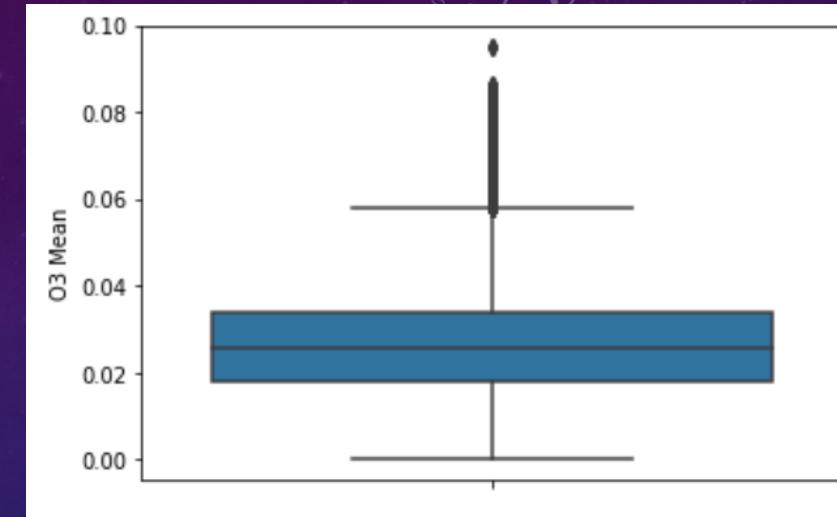
```
##Look Into each value
data=data.dropna()
data.to_csv("/Users/xiangliu/Desktop/CSC560 Data/pollution_AQI.csv",index=True,sep=',')
```

```
data = pd.read_csv('/Users/xiangliu/Desktop/CSC560 Data/pollution_AQI.csv')
data.shape
(436876, 30)
```

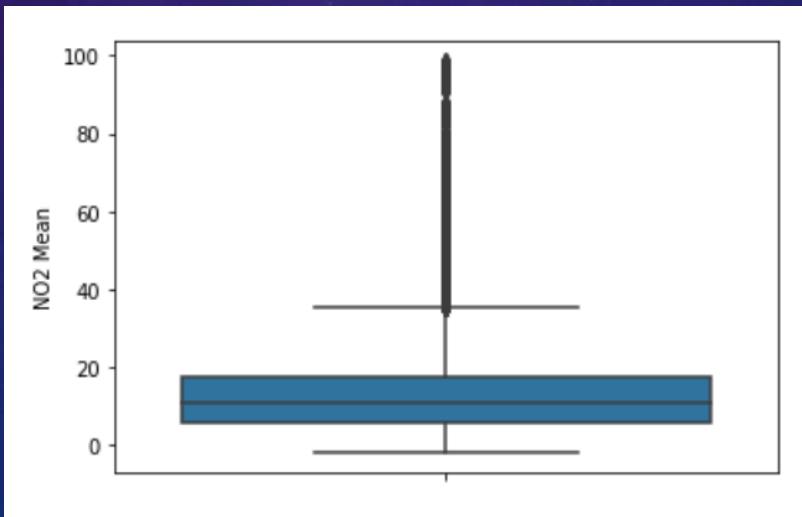
Step 3 Check unusual value



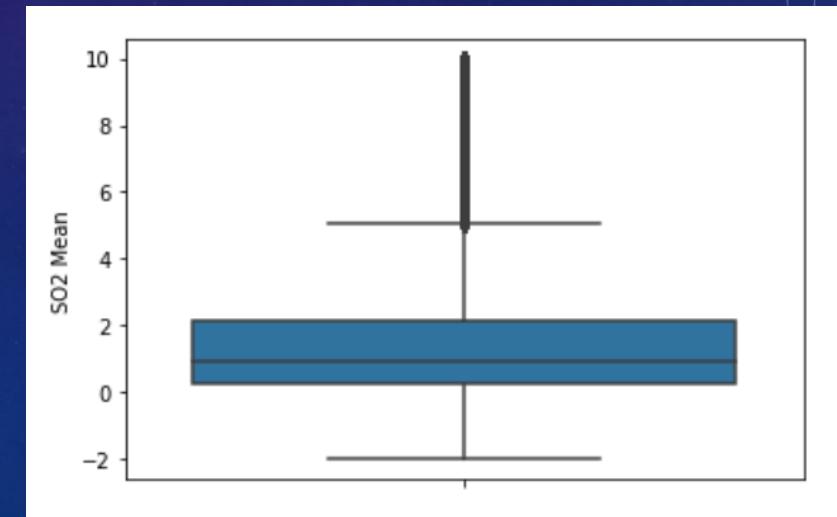
CO Mean



O3 Mean



NO2 Mean



SO2 Mean

HADOOP M/R + HIVE

```
#!/usr/bin/env python
import sys

# input comes from STDIN (standard input)

# input comes from STDIN (standard input)
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # split the line into words
    words = line.split(' ')
    # write the results to STDOUT (standard output);
    # pick the column index 13, 18, 19 for N02 1st Max Hour,03 1st Max Hour,
    #N02 1st Max Hour, C0 1st Max Hour
    print('%d\t%d\t%d\t%d' % (words[13],words[18],words[23], words[28]))
```

```
hadoop jar /usr/lib/hadoop-0.20-
mapreduce/contrib/streaming/hadoop-streaming-2.6.0-mr1-
cdh5.13.0.jar
-input /tmp/data/pollu1.csv
-output /tmp/data/out_put01
-mapper /home/cloudera/myPractice/mapper.py
-reducer /home/cloudera/myPractice/reducer.py
```

```
#!/usr/bin/env python
from operator import itemgetter
import sys

N02=[]
O3=[]
S02=[]
C0=[]
dic_N02={}
dic_O3={}
dic_S02={}
dic_C0={}
current_count = 0
word = None

# input comes from STDIN
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # parse the input we got from mapper.py
    hour = line.split('\t')

    # convert hour (currently a string) to int
    try:
        hour = int(hour)
    except ValueError:
        # hour was not a number, so silently ignore/discard this line
        continue
    # put the four number into four different lists
    N02.append(hour[0])
    O3.append(hour[1])
    S02.append(hour[2])
    C0.append(hour[3])

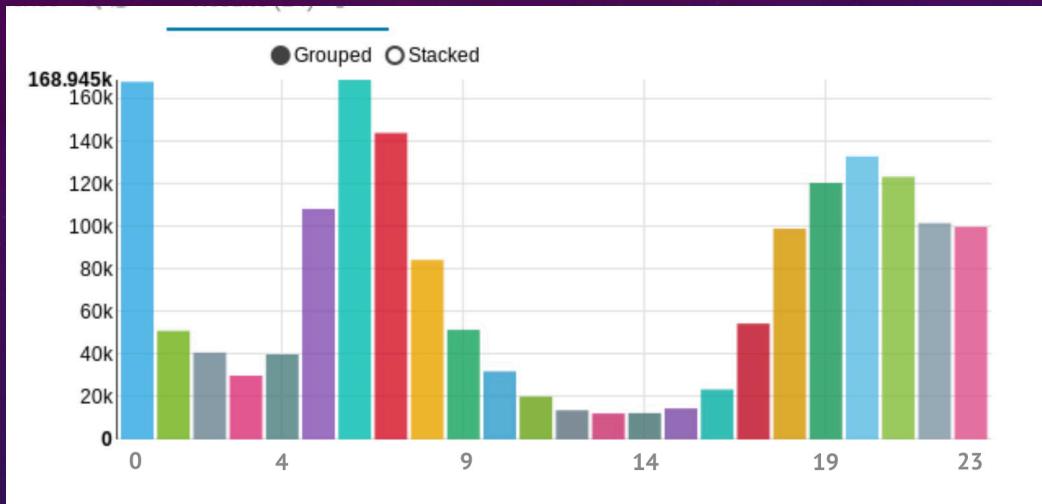
    # read and count the frequency of hours
    for key in N02:
        dic_N02[key] = dic_N02.get(key, 0) + 1
    print(dic_N02[key])

    for key in O3:
        dic_O3[key] = dic_O3.get(key, 0) + 1
    print(dic_O3[key])

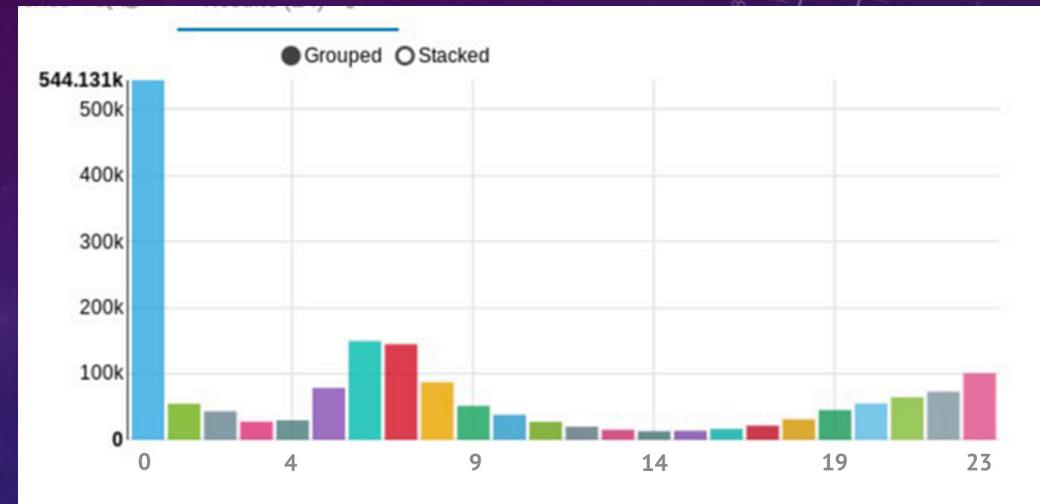
    for key in S02:
        dic_S02[key] = dic_S02.get(key, 0) + 1
    print(dic_S02[key])

    for key in C0:
        dic_C0[key] = dic_C0.get(key, 0) + 1
    print(dic_C0[key])
```

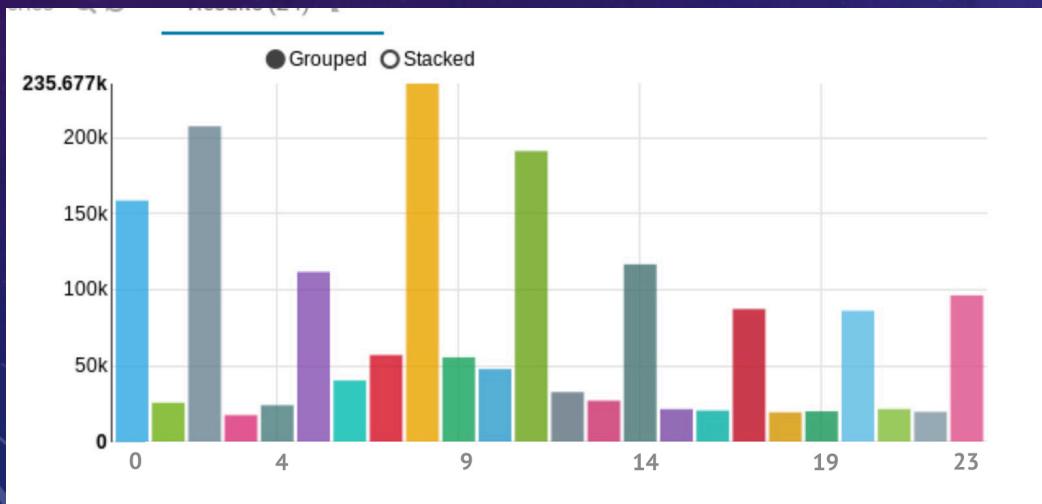
DATA ANALYSIS IN HIVE



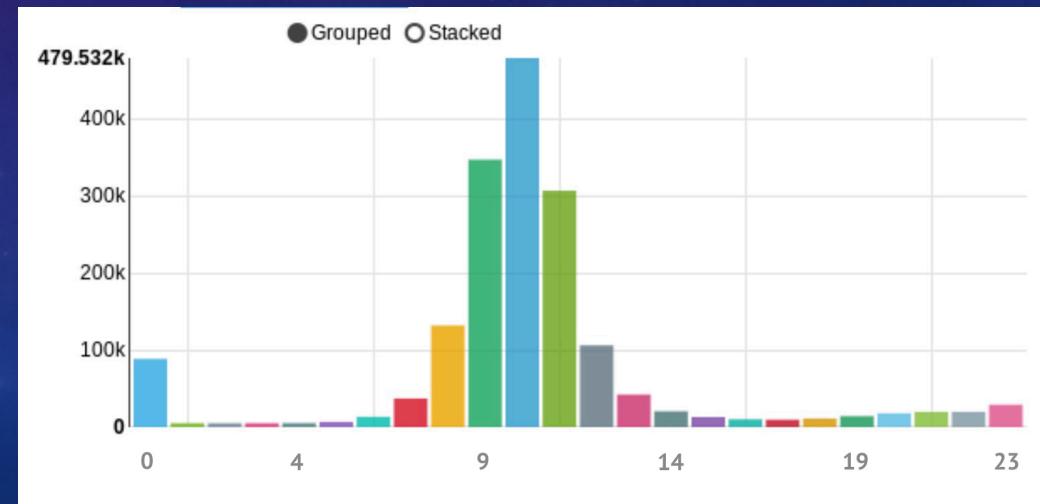
NO₂_1st_max_hour



CO_1st_max_hour



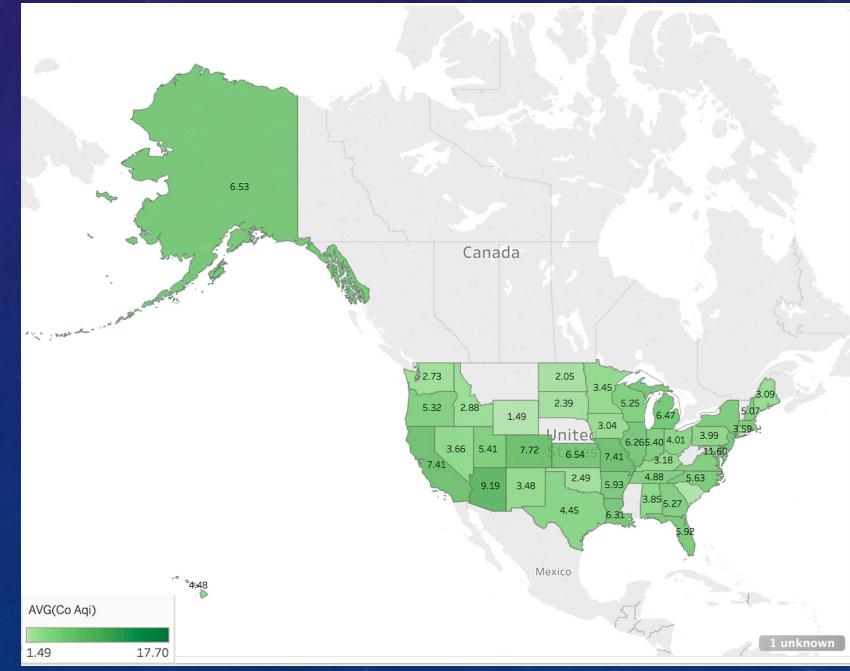
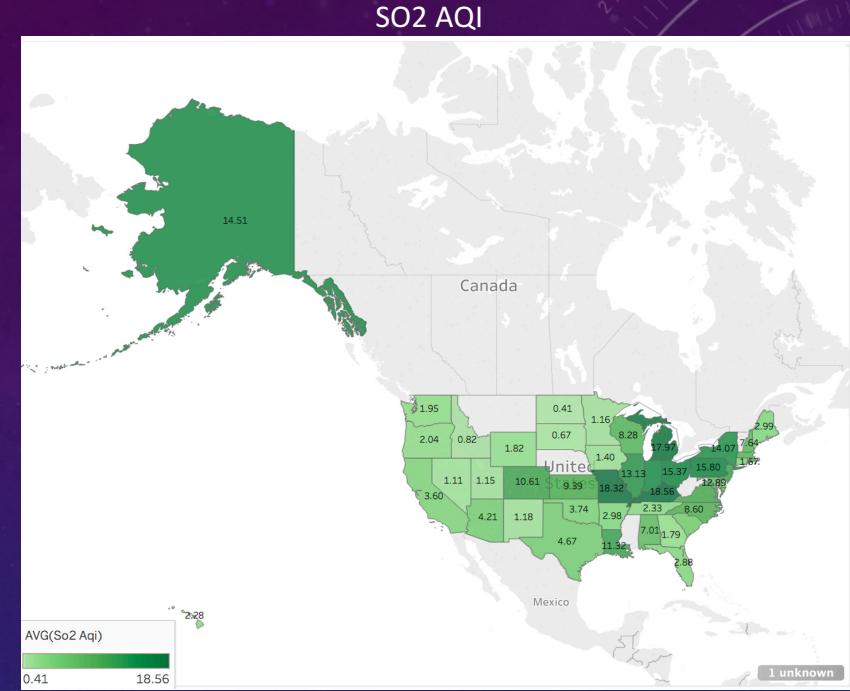
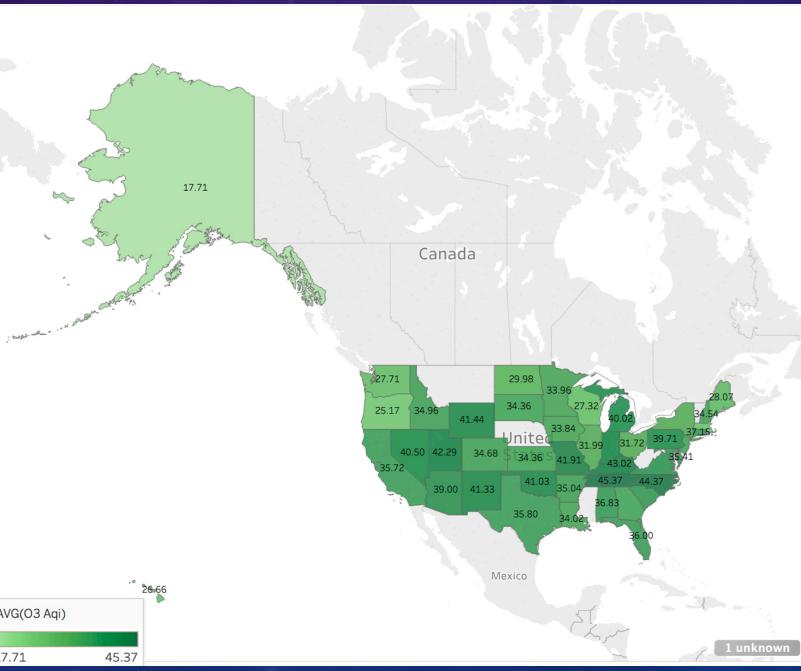
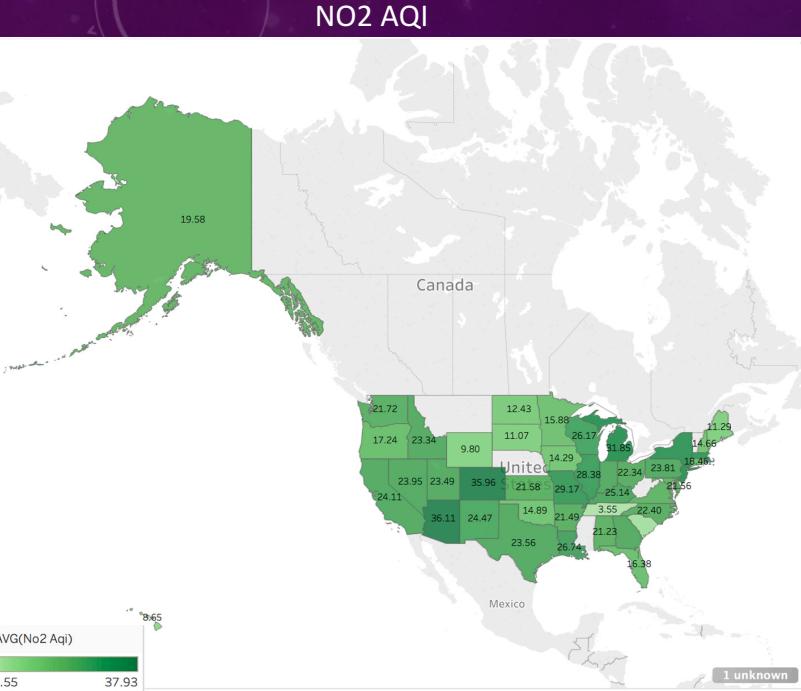
SO₂_1st_max_hour



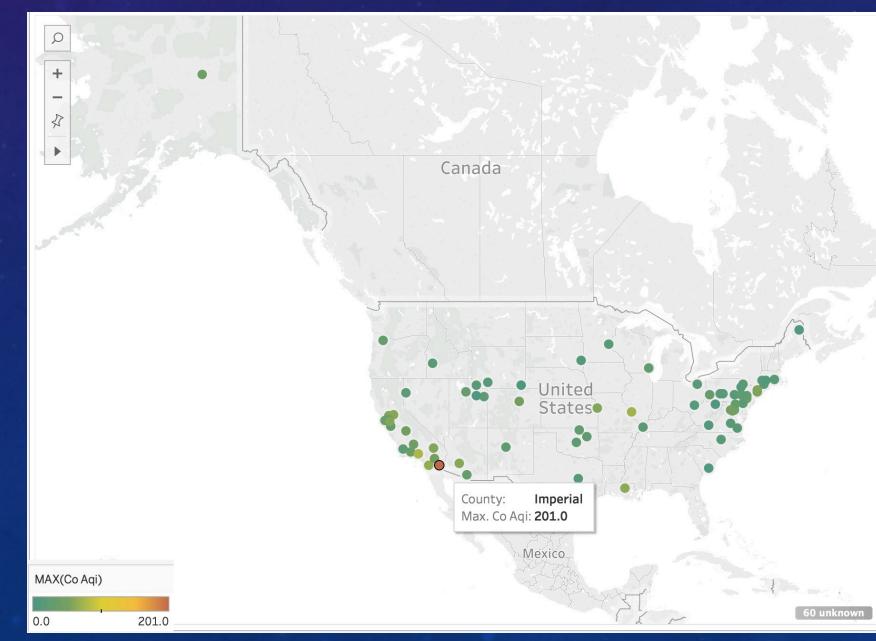
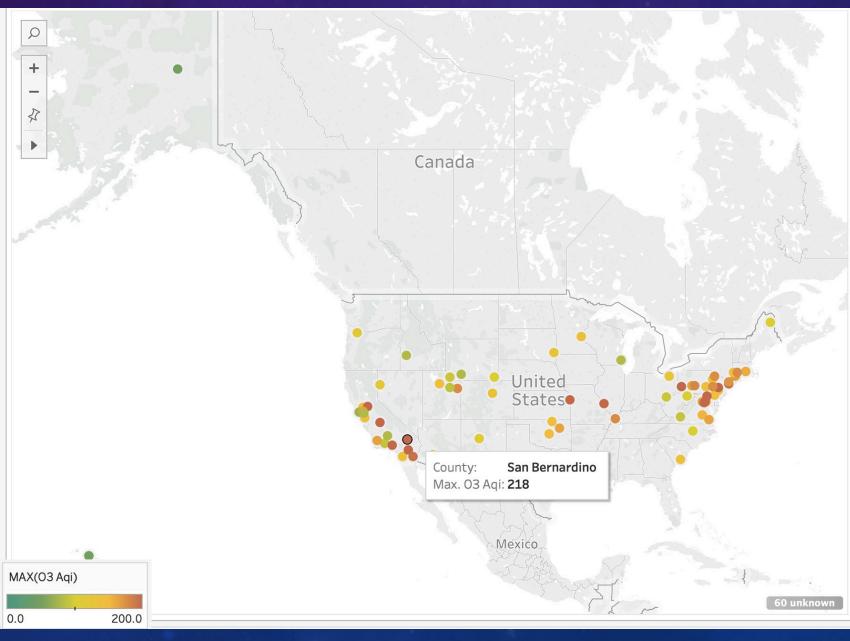
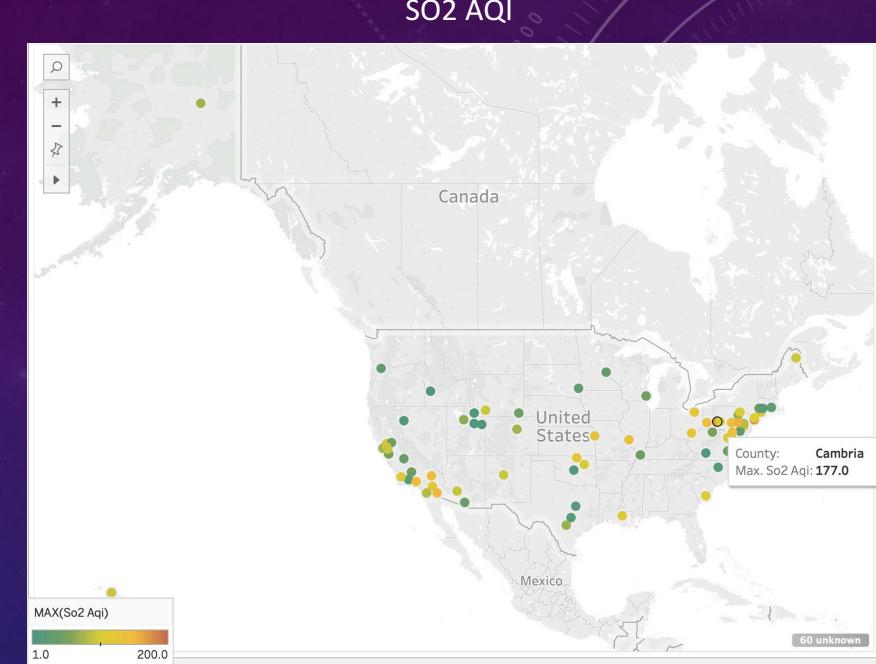
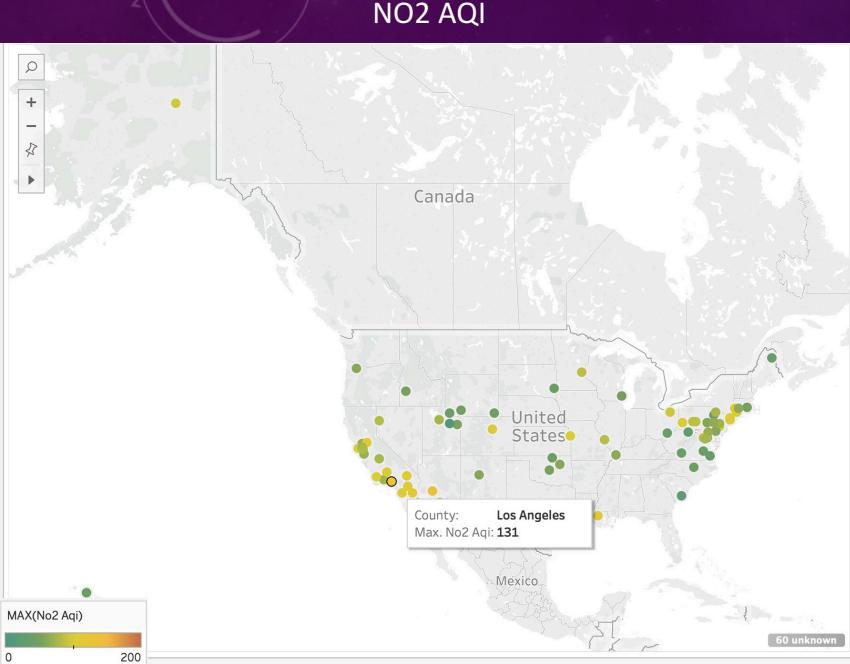
O₃_1st_max_hour

AQI AVG VALUE (STATE)

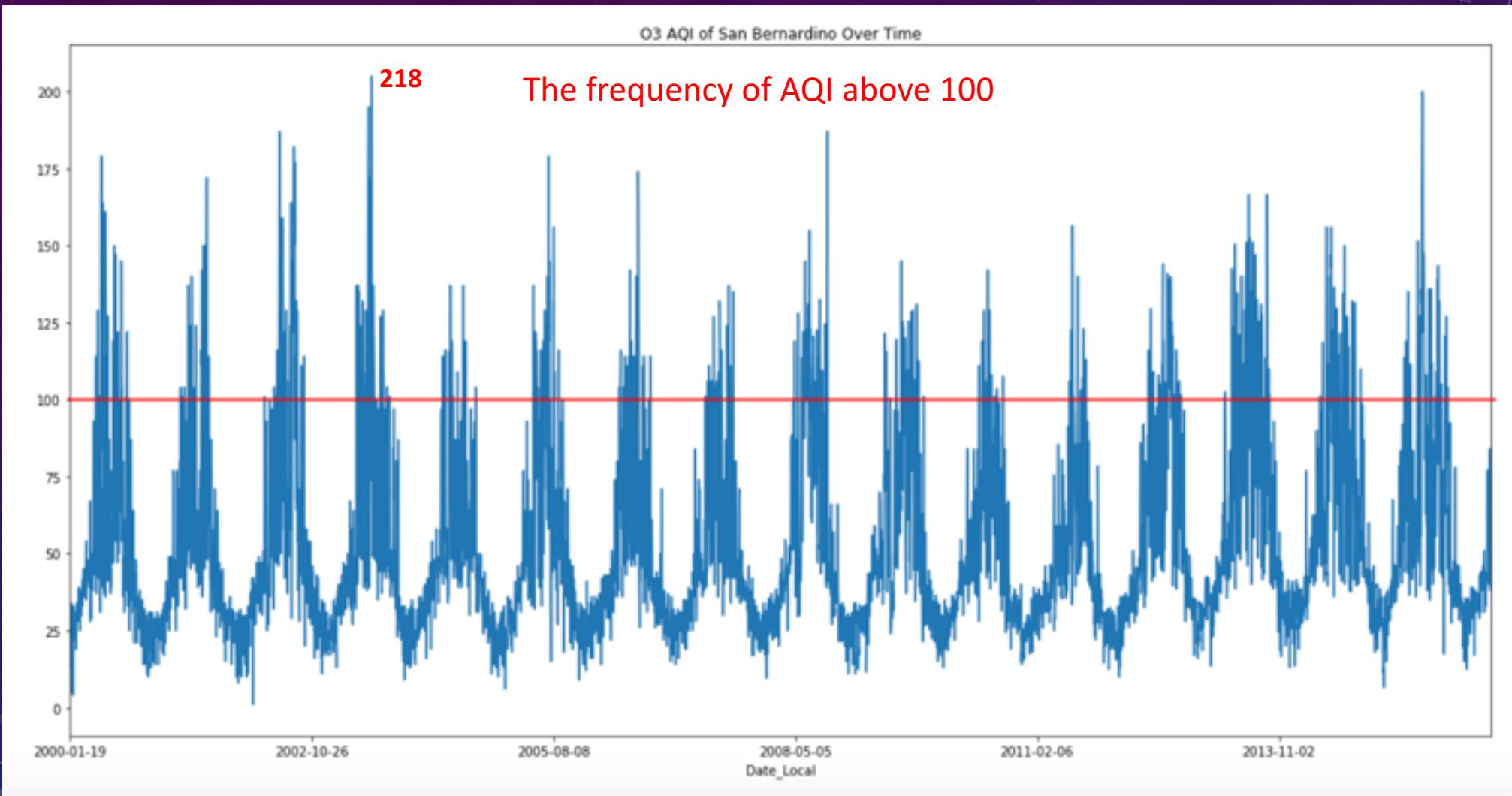
Except Montana, Nebraska, Mississippi,
West Virginia, Vermont these 5 states



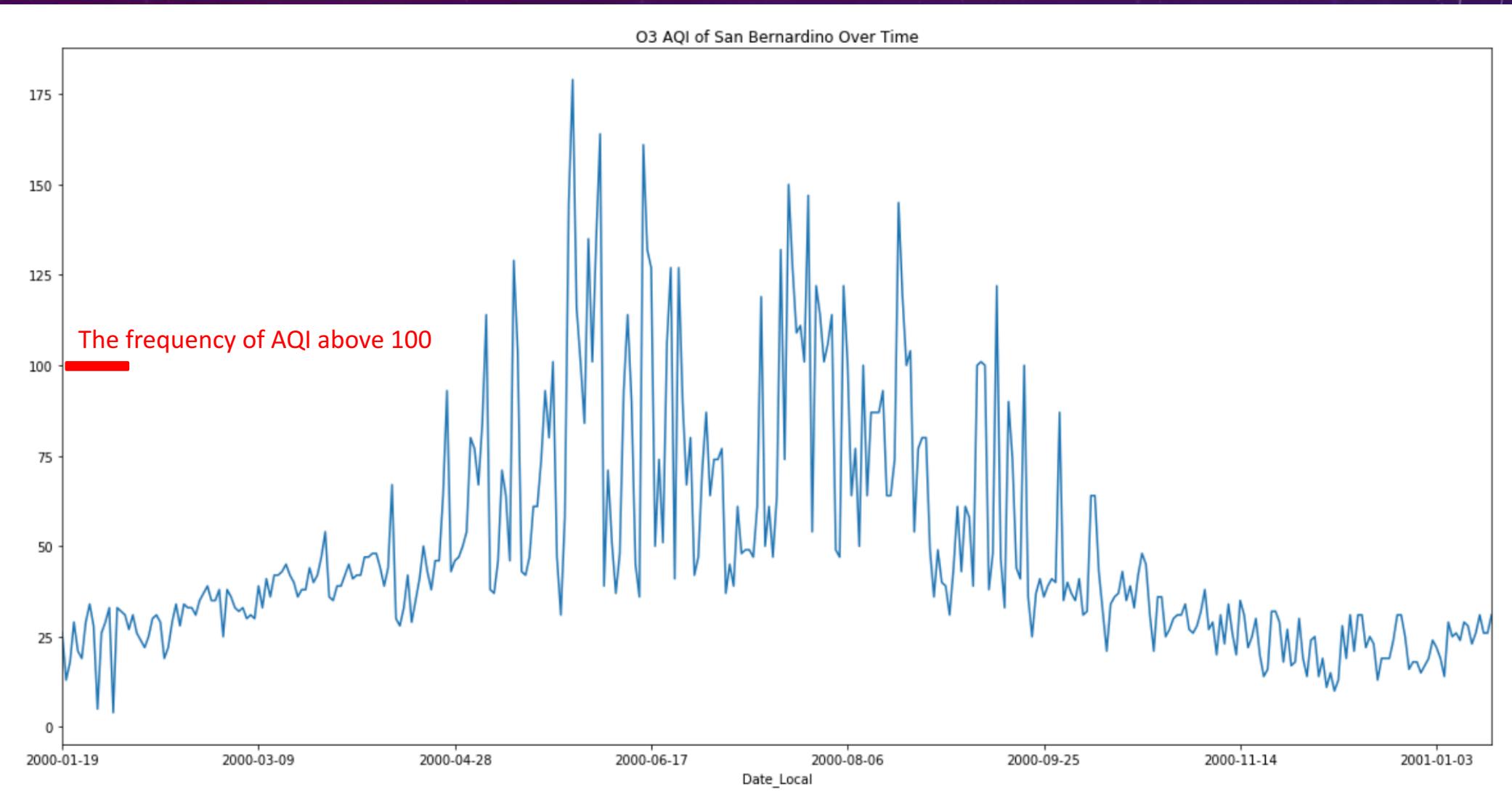
AQI MAX VALUE (COUNTY)



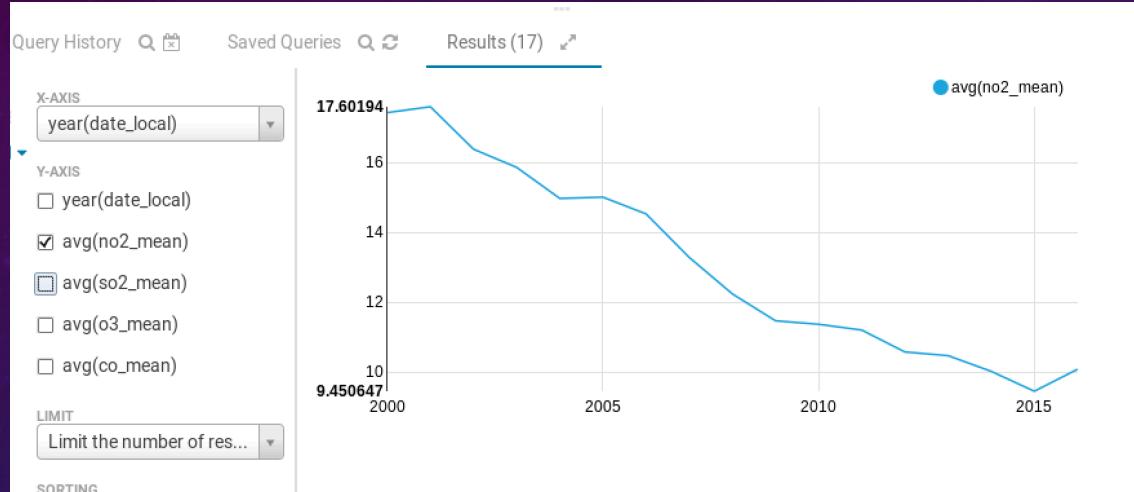
San Bernardino O3 AQI Analysis (2000-2016)



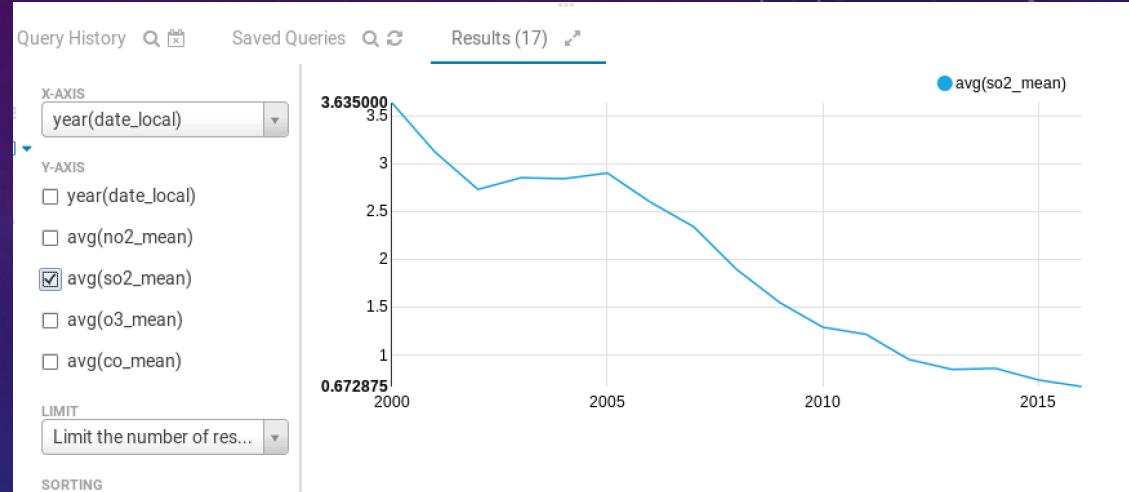
San Bernardino O3 AQI Analysis (2000)



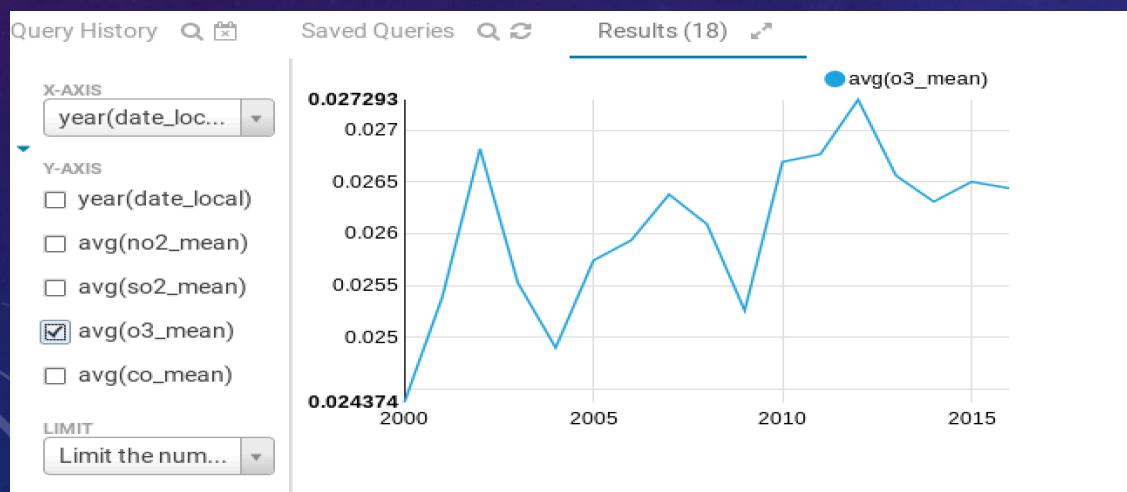
DATA ANALYSIS IN HIVE



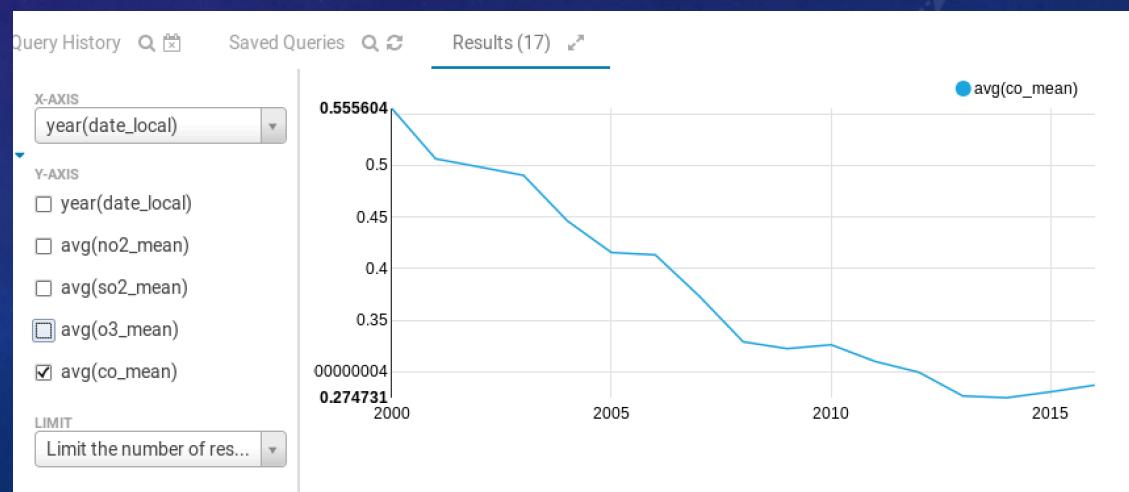
NO2_mean



SO2_mean



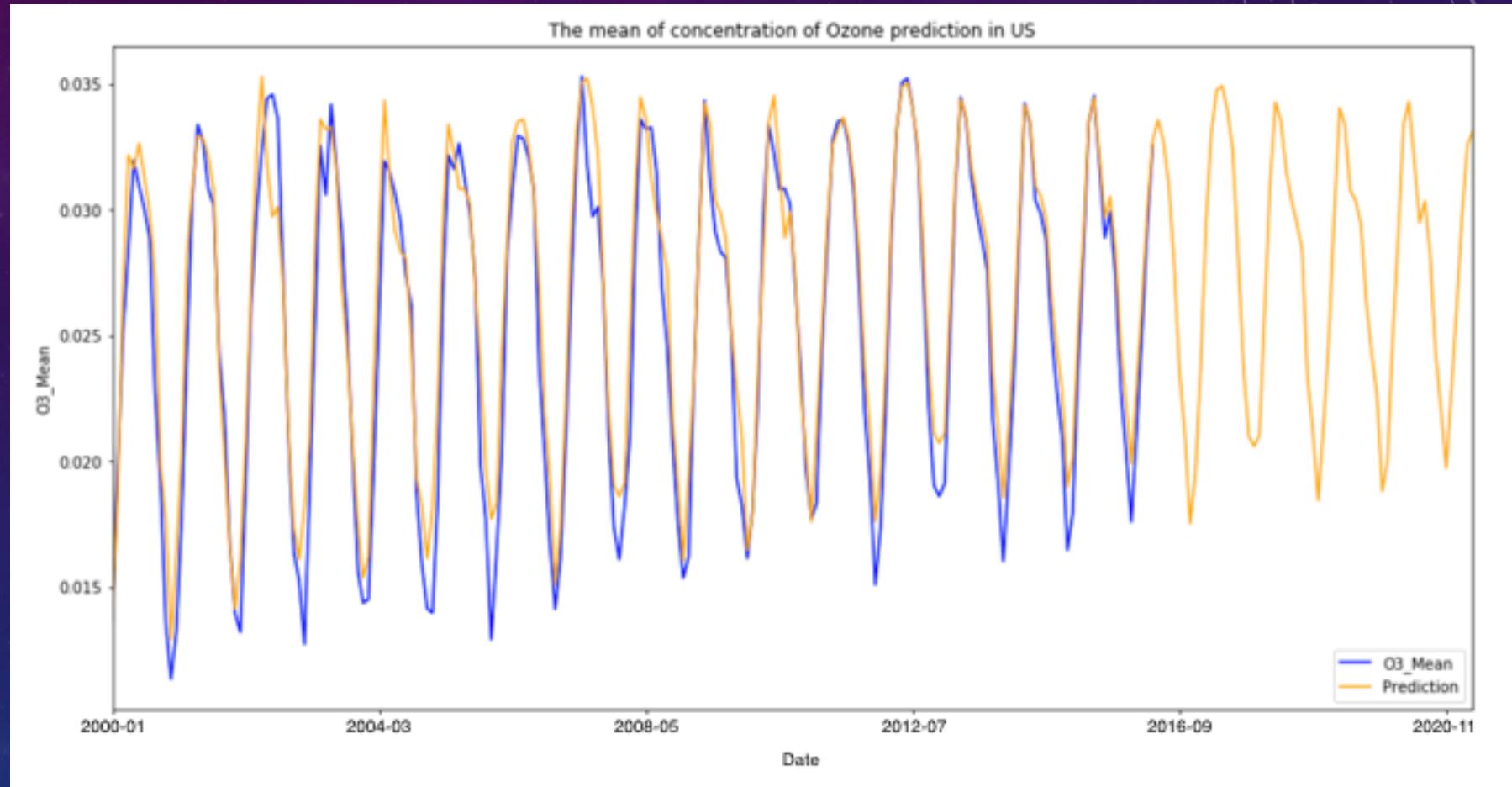
O3_mean



CO_mean

LINEAR REGRESSION MODEL

- Sklearn for the regression algorithm
- Split our testing and training data sets, test_size equal to 20% of the data
- Model accuracy is 0.95



US O3_Mean prediction

SUMMARY

1. Max Value Hour's study

- For NO₂, and CO, most of the places have their max value between 5-9 am and 5- 12 pm
- For O₃, the max value generally appears at 8 am- 1pm
- For SO₂, the peak appears every three hours

2. Mean Value by year study

- NO₂, SO₂ and CO are declining every year gradually
- O₃ is in a stable level with fluctuations

3. Pollution distribution

- Overall, the average AQI for each state is under the safe line (<50)
- Counties in East and West coast have slight pollution on NO₂, SO₂ ; CO is under the safe level; O₃ values in some of the counties have exceed the safe level, might cause health issue

4. Study of O₃ in a specific area (San Bernardino)

- O₃ AQI presents periodical changes, peak appears between May - October

REFERENCE

- https://en.wikipedia.org/wiki/Air_pollution
- <https://enlight.nyc/projects/stock-market-prediction>
- <https://www.kaggle.com/>

THANK YOU FOR YOUR ATTENTION

