



Analysis and Results

Marcella Bona (QMUL), Jia-Chen Hua (QMUL), Adriana Lara (IPN),
Alberto Luviano Juárez (IPN), John Moriarty (QMUL)

IPN Graduate students: Carlos Jiménez González,
Fernando Moreno-Gómez, Royce Richmond Ramírez Morales,
Natan Ismael Vilchis-Tavera

GRADnet Machine learning and AI workshop (2022), online

12/01/2022

Motivations

Pollution has a devastating effect in our lives if we live in big cities.

The threshold for triggering alerts in Mexico City change every four years

- more relaxed than the current acceptable levels set by the WHO
- after 264 days of 2021, the city had reported just the 65% of days as "clean"

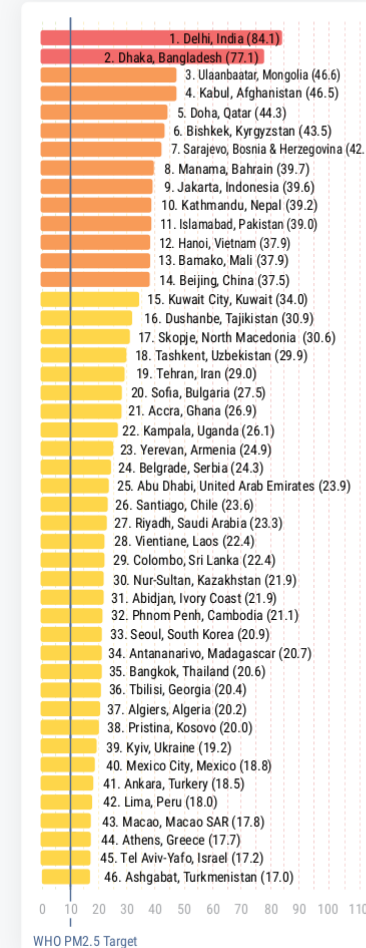
1. White, P. A., Gelfand, A. E., Rodrigues, E. R., & Tzintzun, G. (2019). *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 182(3), 101060.

2. Camahi, Elias, El aire de la ciudad de México supera con creces los límites que la OMS considera peligrosos para la salud, El País, 09/22/2021

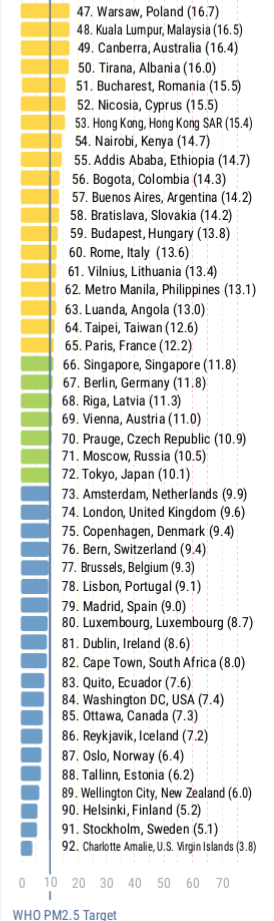
3. www.iqair.com/Fworld-most-polluted-cities%2Fworld-air-quality-report-2020-en.pdf

World capital city ranking

Arranged by annual average PM2.5 concentration ($\mu\text{g}/\text{m}^3$)



World air quality report 2020 from QAir



Motivations

“Smart cities” can gather an enormous amount of data to help them improve the situation

Big data capabilities being built in many contexts allow to store, clean, and analyse these data.

Data can be analysed and with statistical modelling and machine learning we can learn behaviours and obtain predictions

“Smart cities” can imagine solutions and empower society to act, citizens and governments

Most data sets can come directly from the citizens themselves: traffic for example

Basic data: traffic data from google is free to download (up to some level)

SAPIENS: can we use the basic free level of traffic data to learn about pollution?

Data

Pollution data: CDMX Data Agency provides pollution data in terms of:

- 27 stations, levels of 9 pollutants recorded every hour
- Data Agency provides a clean set of data after 3 months of being taken

Traffic data: basic google images which are free and available instantaneously

- Downloaded in the SAPIENS database 3 times per hour.
- Traffic information google images of the 10 km² map around each of the sensors are downloaded.

Data processing

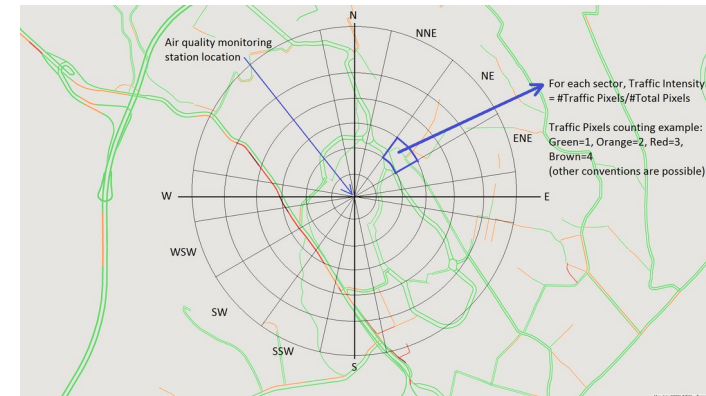
Pollution data: need an extra clean-up to remove outliers and to check the number of effective measurements for each sensor/station:

- Lots of sensors have very few data
- Identified 3 sensors with high number of measurements for each pollutant
- The data from these 3 sensors are used in our modelling analysis
- The other sensors with a minimum of measurements to be used for validation

Traffic data: images are translated into traffic intensity measurements based on concentric circles around the sensor position

- Information in terms of thickness of traffic colour-labelled lines or line segments: e.g. a wider street labelled orange is likely to have more traffic than a narrower street labeled by the same traffic colour
- Count pixels of traffic colors to quantify traffic flow or volume.
- take other non-traffic pixels into consideration: if a station is located far from streets, the traffic intensity surrounding should be lower as those non-traffic pixels are not producing or “emitting” pollutants.

id_station_id	Null all Day	Null of 6 to 20 hrs
MER	489	170
CAM	962	485
PED	1541	779
IMP	2136	1246
TLA	2586	1391
ARA	4272	2492
LVI	4272	2492
VAL	4272	2492
SAG	7901	4607
SFE	10246	5990



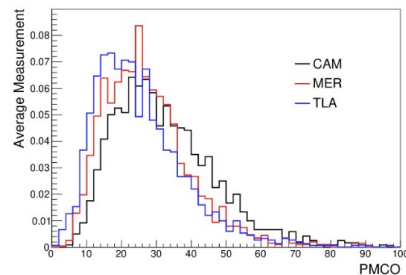
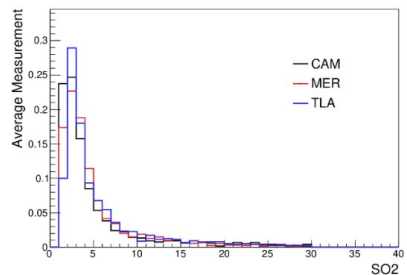
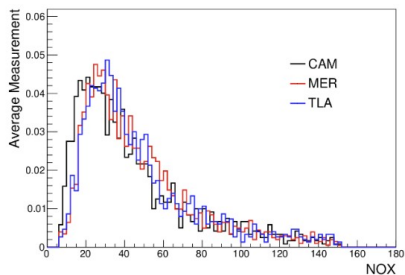
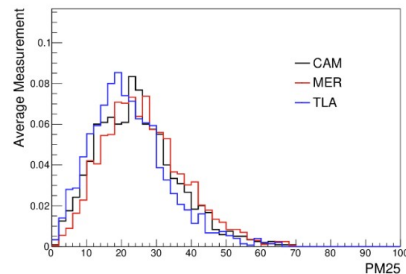
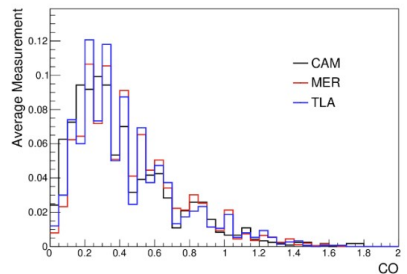
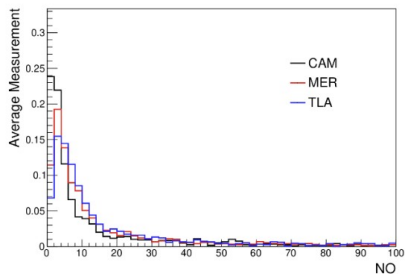
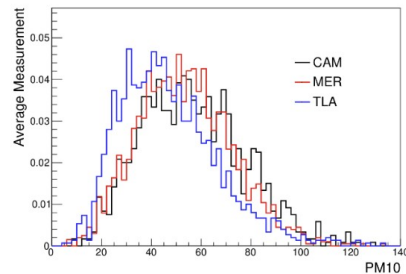
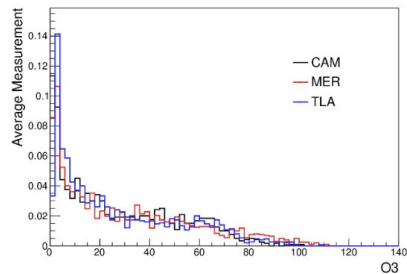
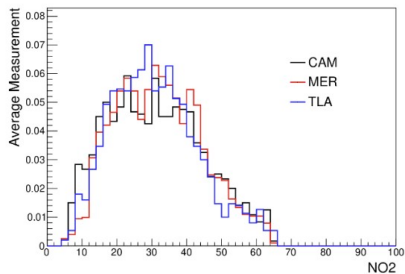
Pollutant Data Analysis

List of pollutants:

- PM10: particulate matter with 10 μm or less in aerodynamic diameter ($\mu\text{g}/\text{m}^3$)
- PM2.5: particulate matter with 2.5 μm or less in aerodynamic diameter ($\mu\text{g}/\text{m}^3$)
- PMCO: particulate matter with aerodynamic diameters between 2.5 and 10 μm ($\mu\text{g}/\text{m}^3$)
- SO₂: Sulfur Dioxide (ppb)
- O₃: Ozone (ppb)
- CO: Carbon Monoxide (ppm)
- NO₂: Nitrogen Dioxide (ppb)
- NO: Nitrogen Monoxide (ppb)
- NO_x: Nitrogen Oxides (ppb)

$\mu\text{g}/\text{m}^3$ = micrograms per cubic metre
ppb = parts per billion
ppm = parts per million

Pollutant Data Analysis

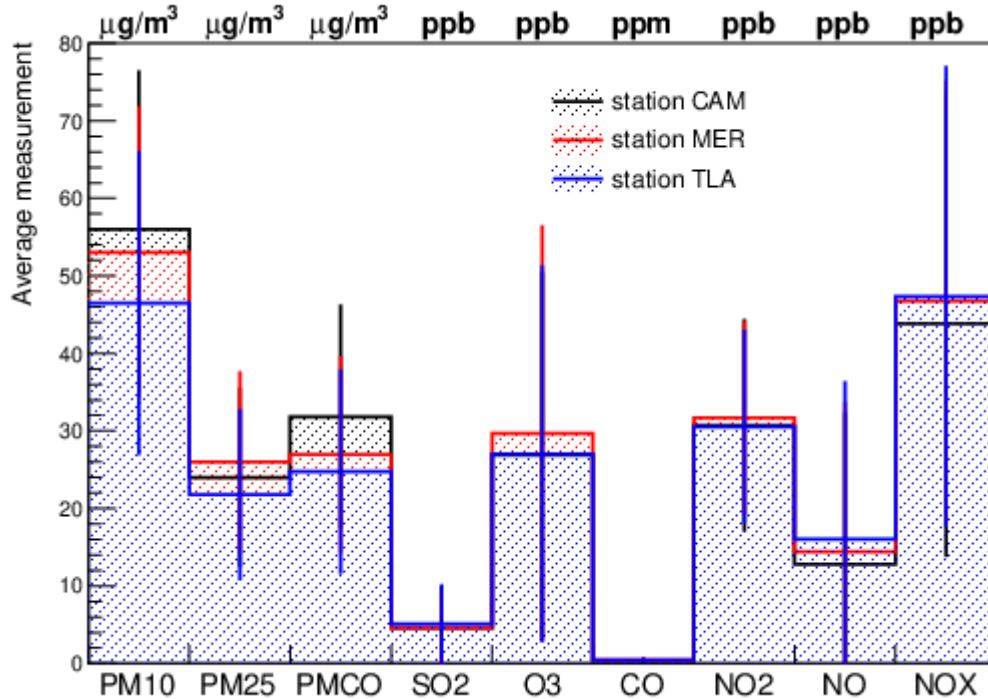


Distributions of the pollutant measurements:

Comparison between the three sensor stations (labelled: CAM, MER and TLA)

Compatible distributions across the stations

Pollutant Data Analysis



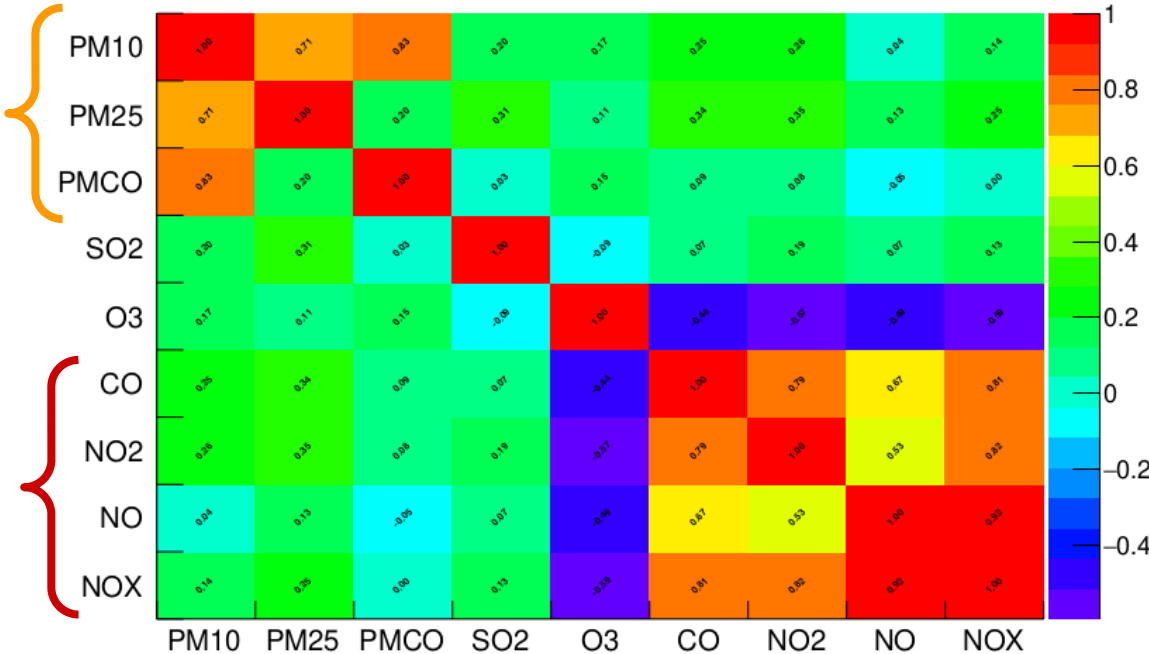
Average measurements for each pollutant and for each station used.

Error bars are the RMS of the pollutant distributions.

Units differ depending on the pollutant considered.

Pollutant Data Analysis

Correlations



Correlation between pollutants:

Red colours: positive corr.

Blue colours: negative corr.

Two groups:

1: {PM10, PM25, PMCO}

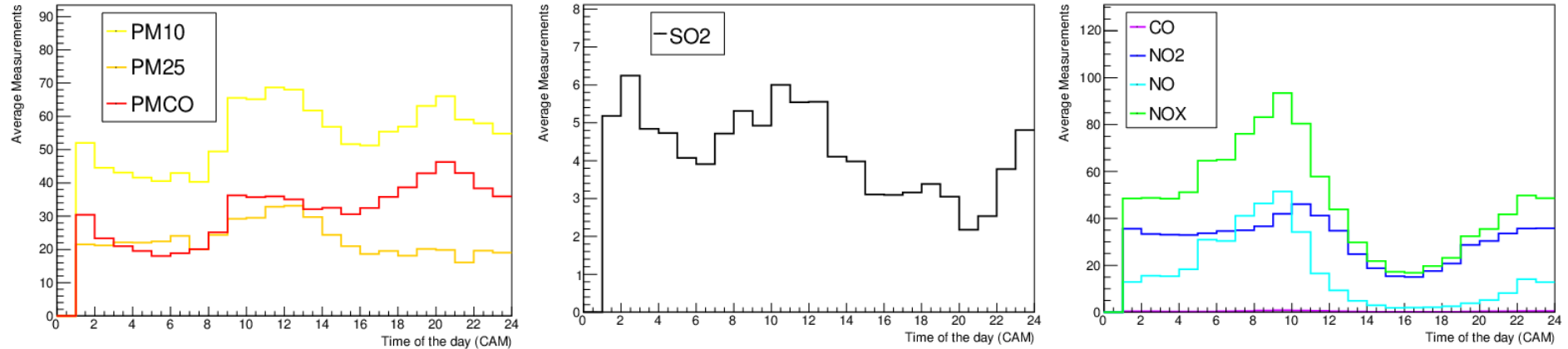
2: {CO, NO2, NO, NOX}

O3 anticorrelated with **group 2**

SO2 uncorrelated with the others

Pollutant Data Analysis

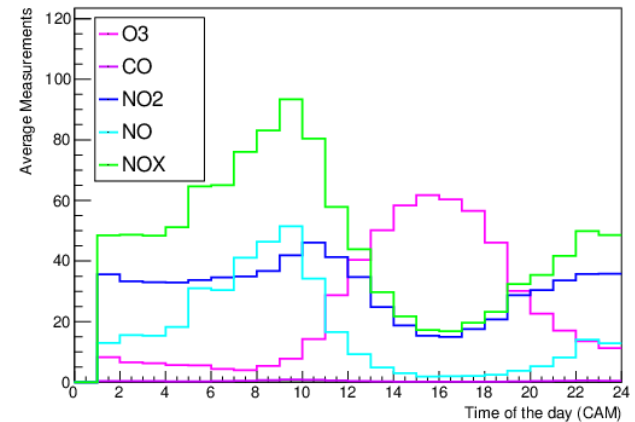
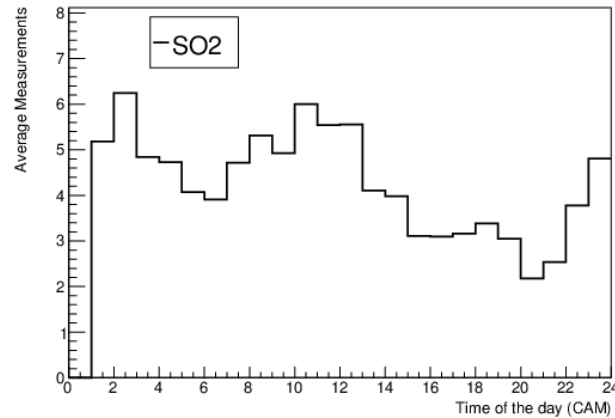
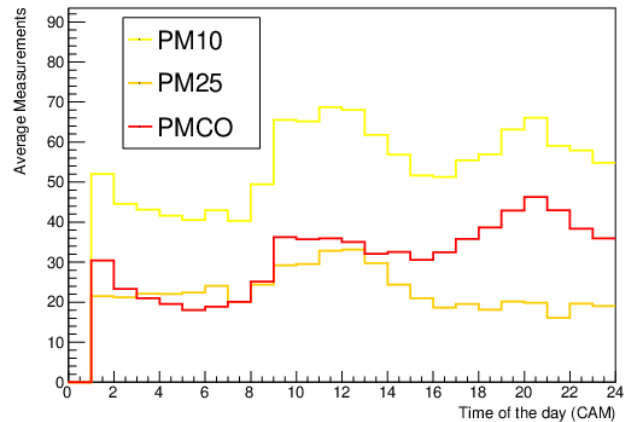
Check modulations during the day in the pollution measurements



Pollutants in correlated groups with their absolute measurements.
Average measurements shown for each hour of the day.

Pollutant Data Analysis

Check modulations during the day in the pollution measurements

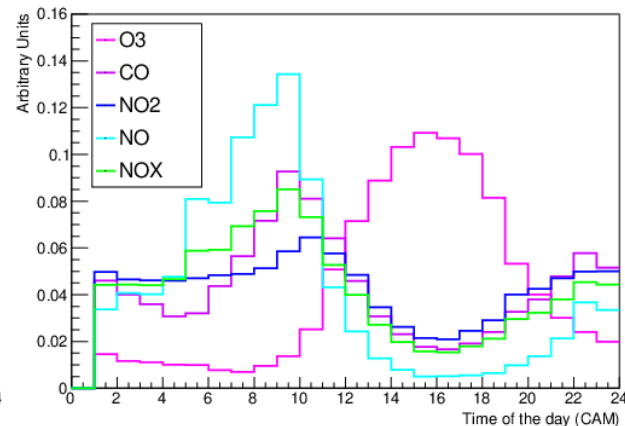
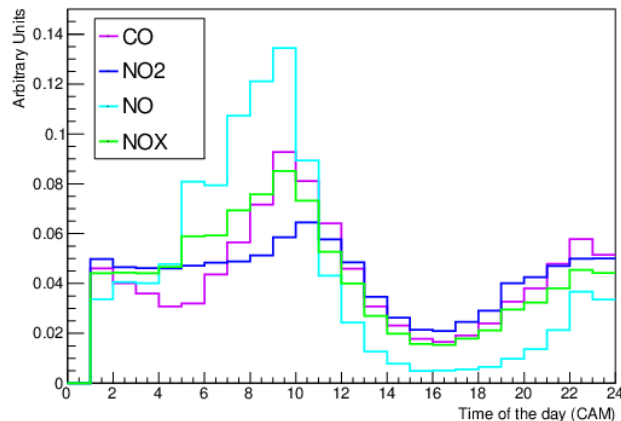
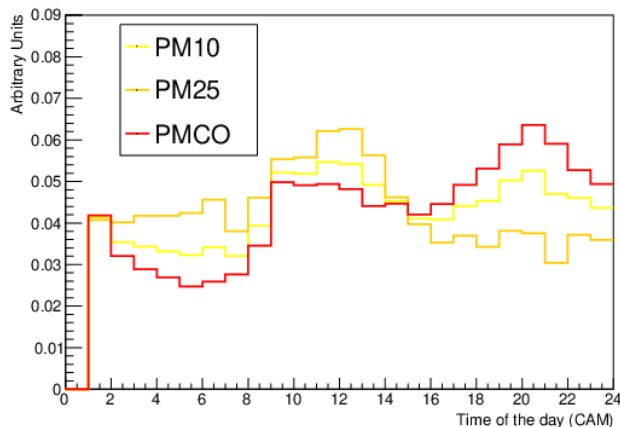


Adding anti-correlated O3

Pollutants in correlated groups with their absolute measurements.
Average measurements shown for each hour of the day.

Pollutant Data Analysis

Check modulations during the day in the pollution measurements



Pollutants in correlated groups:
here normalised distributions for shape comparison.

Traffic Intensity Model

Google Maps images with traffic layer/colouring: green, orange, red, dark red
 % of coloured pixels in annular sectors -> traffic intensity; physical density “field”

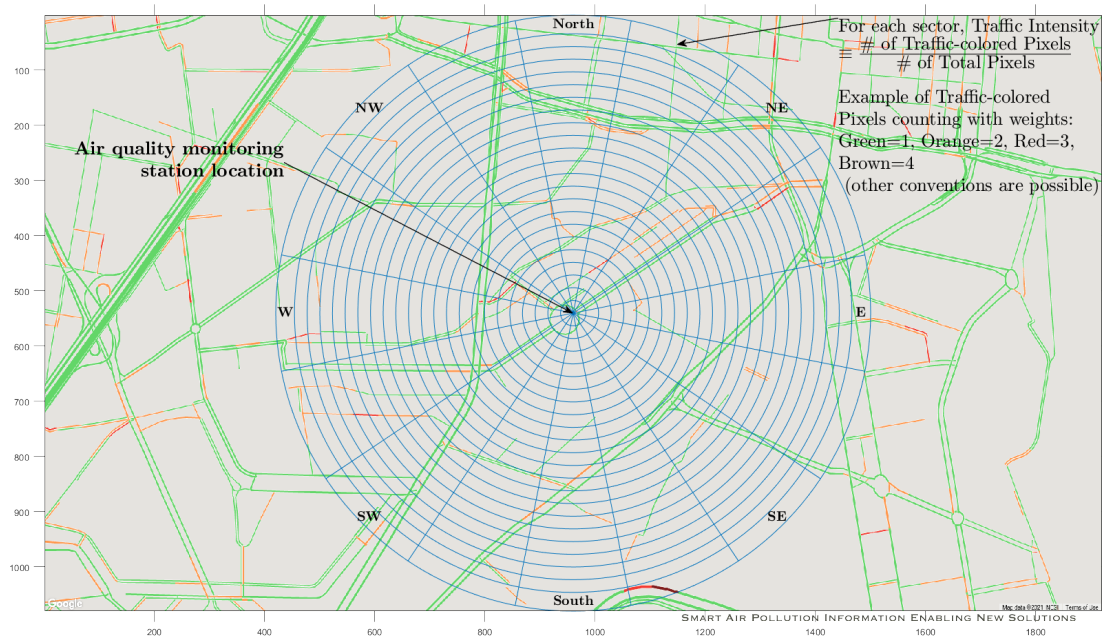
Dimension reduction
 from HD image 1920x1080:

(4 colours)x(16 angles)x(118 rings)

Or aggregating the angles:

(4 colours)x(118 rings)

Only 23 rings in this plot
 118 rings shown in the next page



Traffic Intensity Model



Width of 118
concentric
rings: 10m

Traffic Intensity Model



Details of the mathematical model over to Jia-Chen..

Outlook and Conclusions

SAPIENS has built a database with both pollution measurements and traffic images, so we have:

- Cleaned and analysed the data and identified patterns
- Developed a model to extract the traffic intensities from Google Map images
- Used the regression modeling to (1) obtain interpretable insights on the relation between traffic and pollutants; and (2) train it on the data from three stations (traffic and pollution data) and cross-validated it to avoid overfitting

On-going activities:

- Validation/testing phase: use other sensors data to validate/test model
- Paper in preparation

Outlook and Conclusions

There are more ideas and more possibilities to exploit and learn from these data.

More ideas on how to exploit the predicting power of the modeling

E.g., incorporating meteorological data, going beyond linear modeling techniques, etc.

Stay tuned for more from us

