# Phân tích ô nhiễm không khí ở VN

## *Sử dụng R & Openair package*

**Tuan Vu**

*Senior air quality scientist, King's College London*
*Email: tuan.vu@kcl.ac.uk*
*https://tuanvvu.github.io/*

*10-April-2020*

# Aims

- Introduction to R data analysis software

- Introduction and use of the R package: *openair*

- Machine learning in data analysis

*"We can only see a short distance ahead, but we can see plenty there that needs to be done"- Alan Turning*

# I. Introduction to R data analysis software

## 1. Downloading and installing R/ Rstudio

## 2. General approach to data analysis
➢ Use scripts: save all objects in the current R sessions as an .RData file
➢ Leave the data alone: as much as you can
➢ Coding style
➢ Simple R and vectors: R cheat sheet

## 2. Useful packages
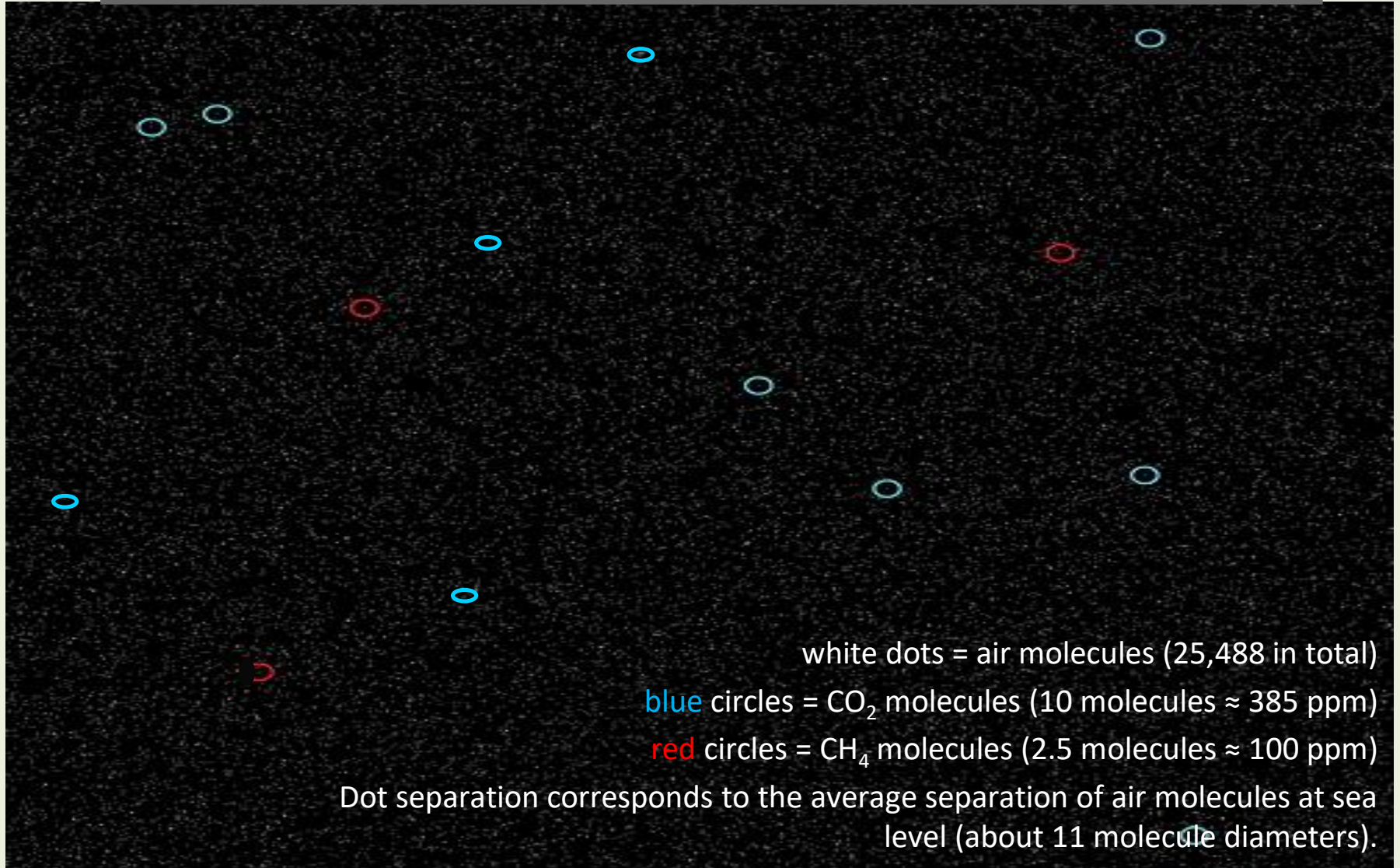➢ lubridate/dplyr/plyr
➢ ggplot2
➢ openair/worldmet

# II. Introduction to *"openair"*
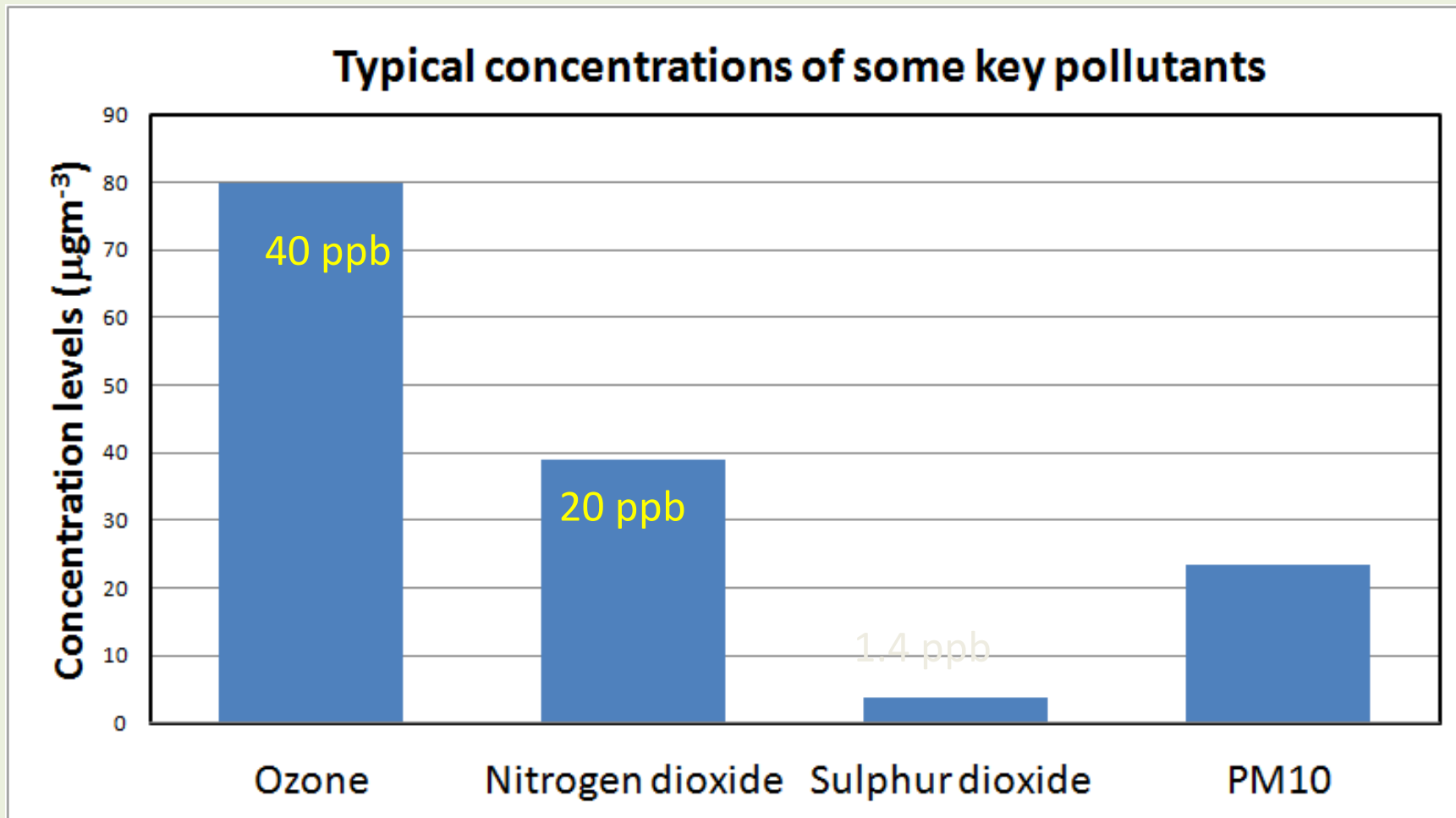
## Useful *openair* functions

1. Summary data: Understand your data

2. Merging data sets

3. Selecting data by date

4.  Averaging data to different time intervals

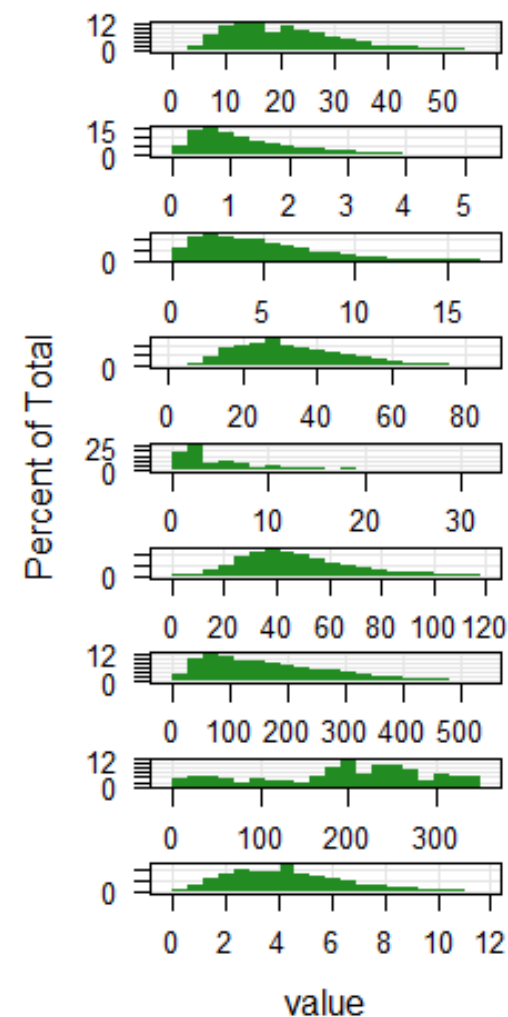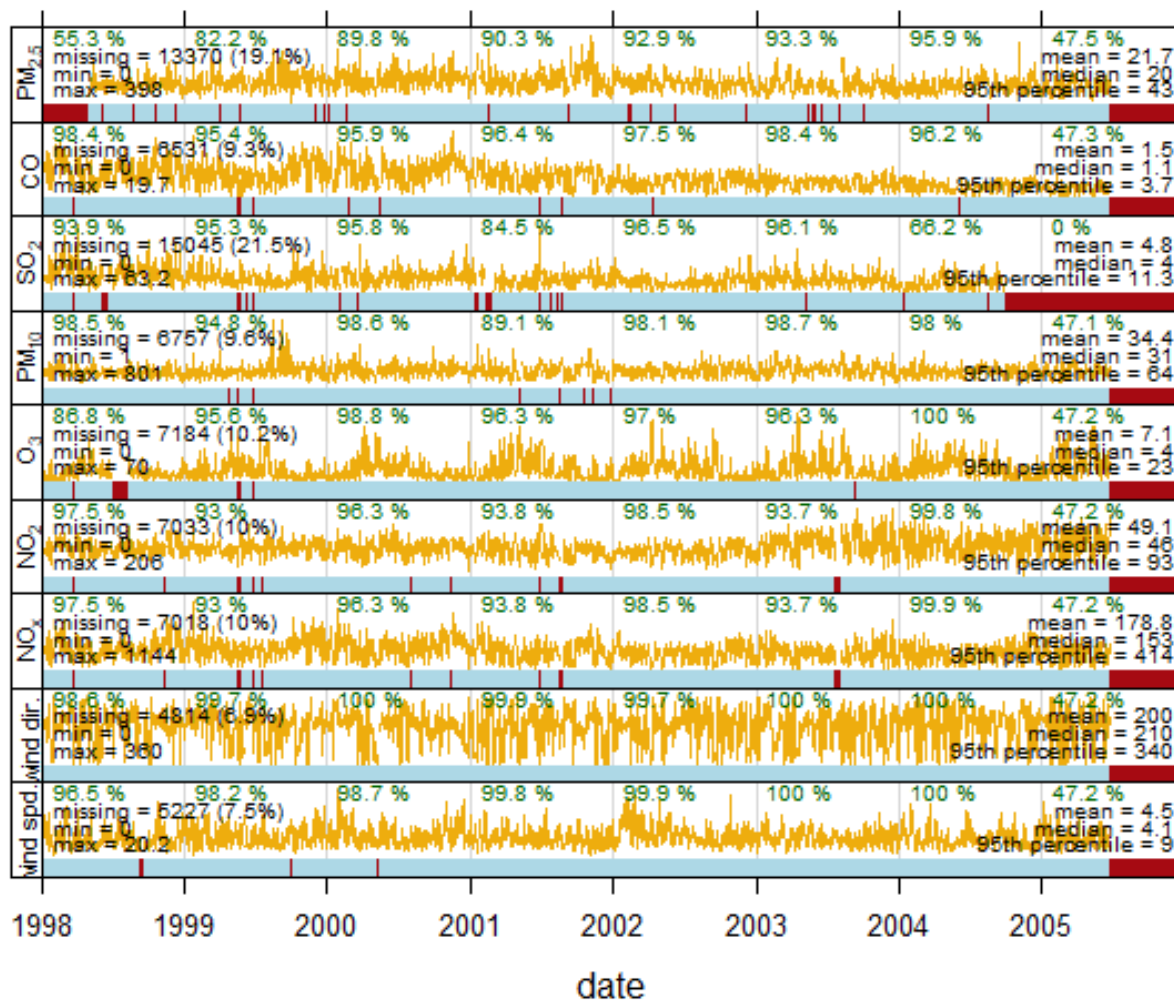5. The *ScatterPlot*

# Understanding about the pollutants



white dots = air molecules (25,488 in total)

blue circles = $CO_2$ molecules (10 molecules ≈ 385 ppm)

red circles = $CH_4$ molecules (2.5 molecules ≈ 100 ppm)

Dot separation corresponds to the average separation of air molecules at sea level (about 11 molecule diameters).

Link to rice comparison

# Levels of Pollutants



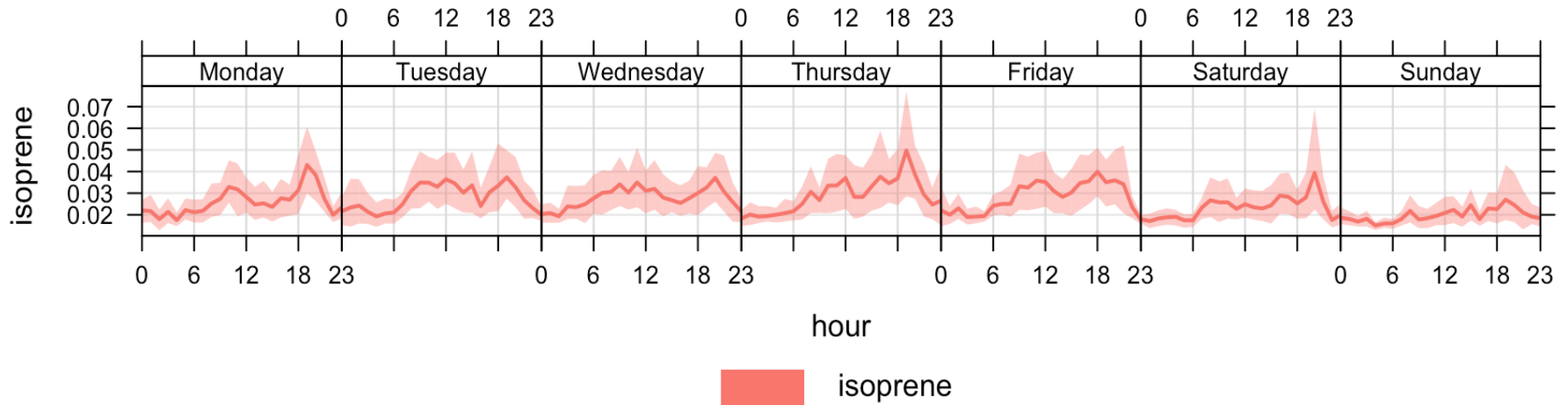**Typical concentrations of some key pollutants**
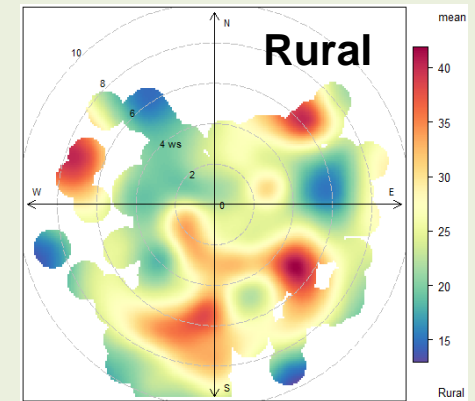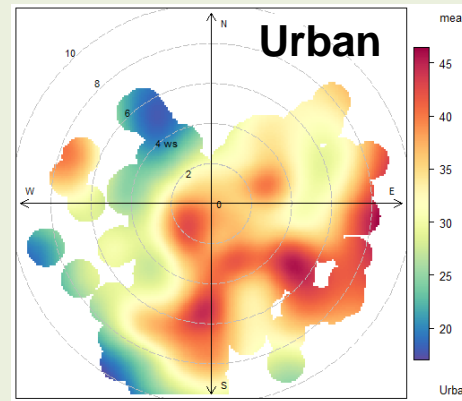
Data obtained from: http://uk-air.defra.gov.uk/
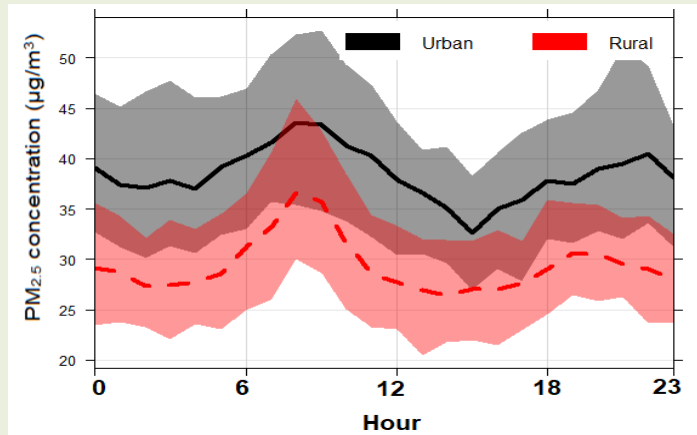
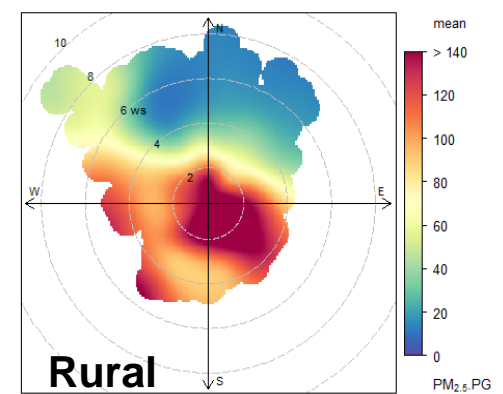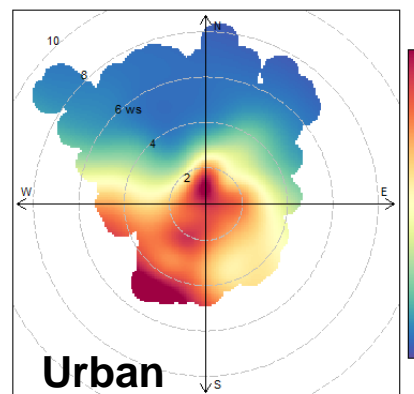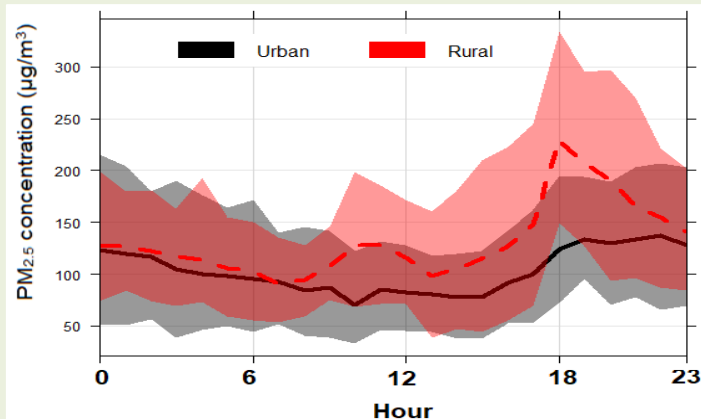# summaryPlot

# 6. CorPlot

# 6. Time Variable Function



mean and 95% confidence interval in mean

# PM$_{2.5}$ diurnal patterns & Polar Plot
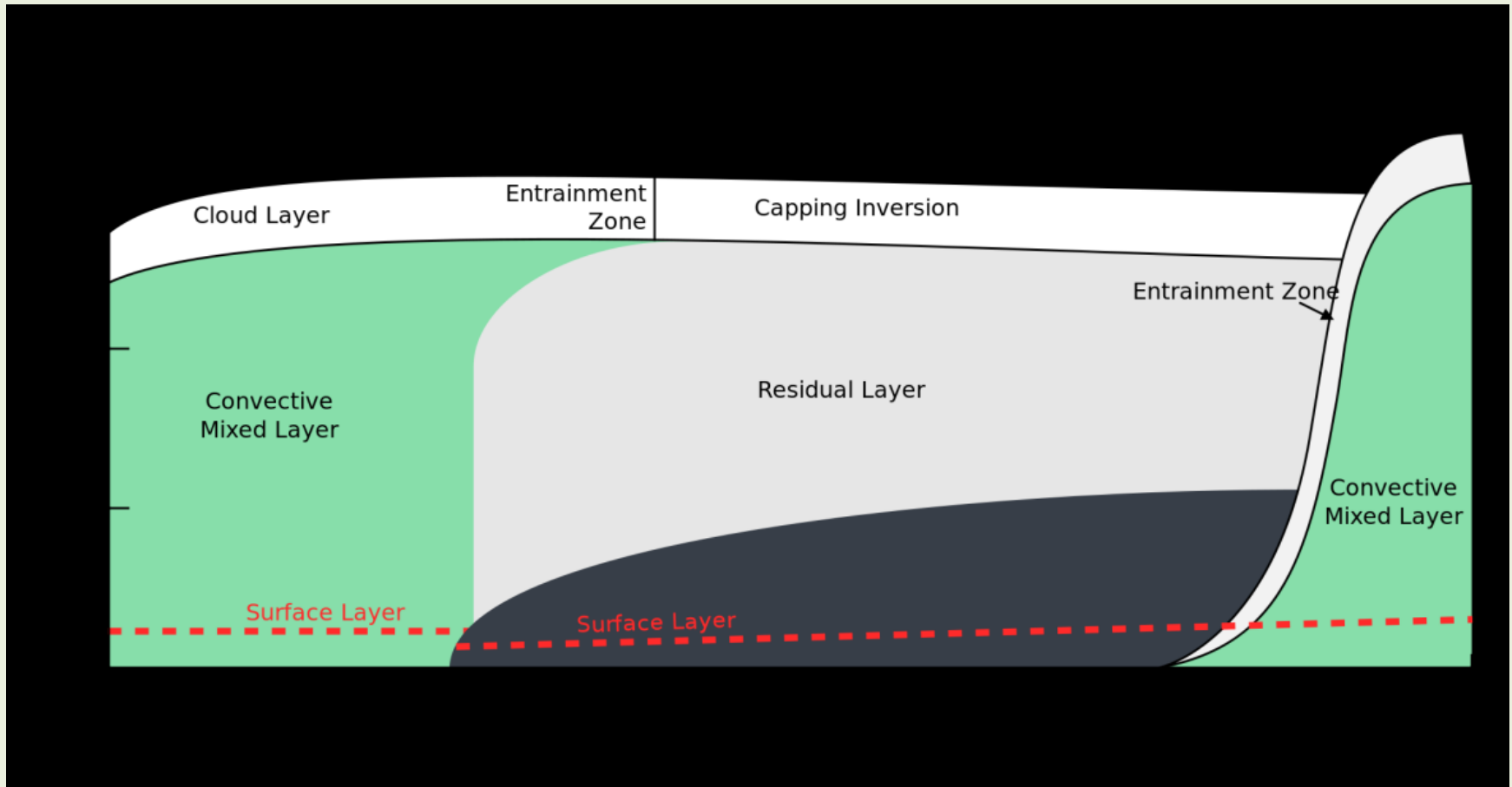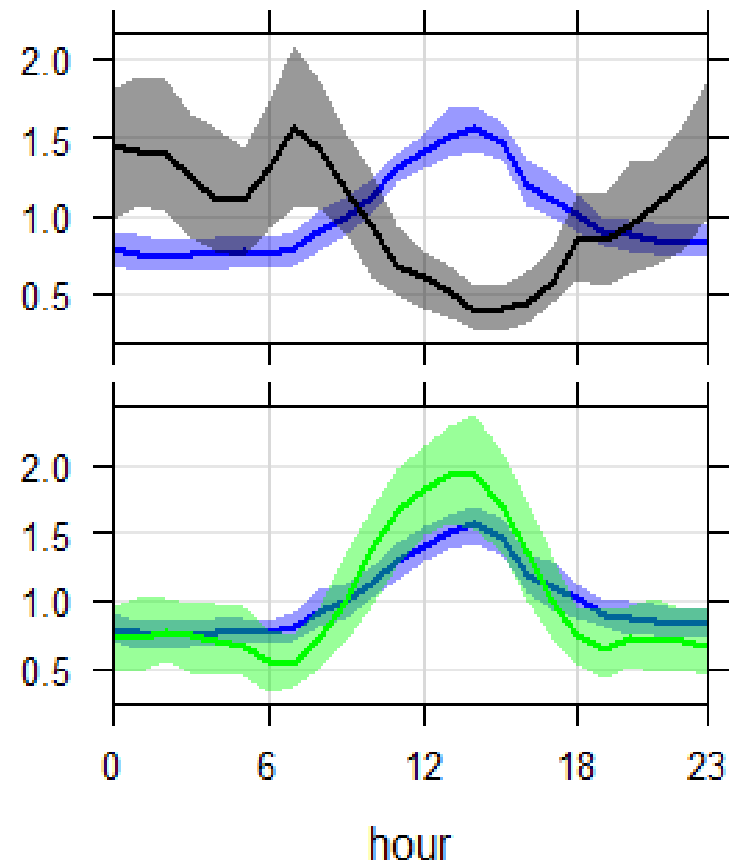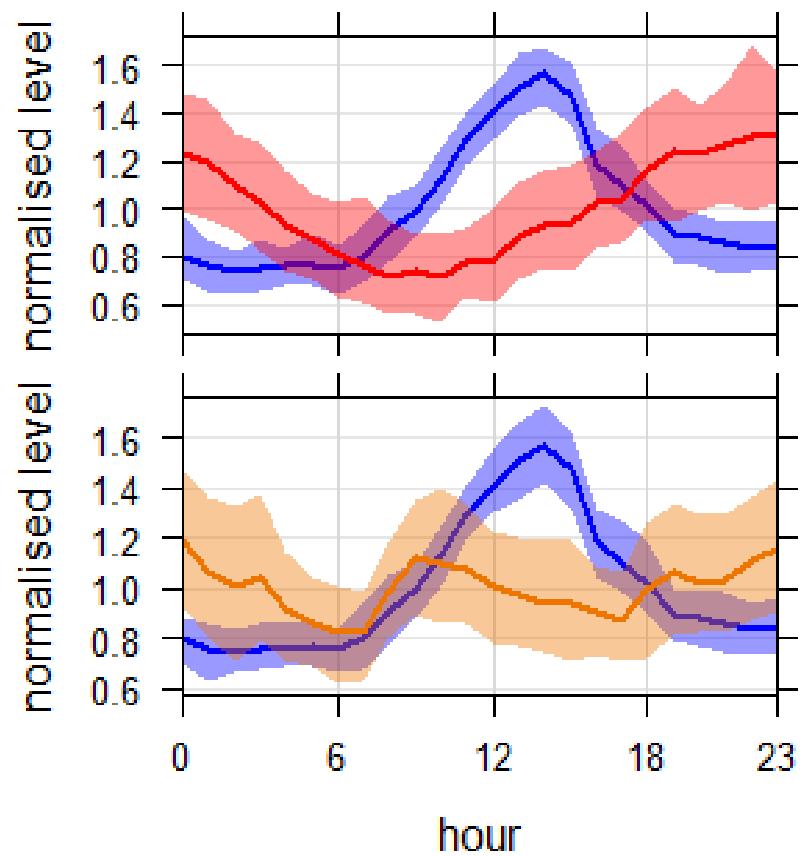
## Summer 2017



## Winter 2016



*Hourly PM2.5 data in Pinggu was provided by PKU*

# Diurnal variations in the ABL



**Cooling of surface during sunset – stable nocturnal layer grows from below**

**Heating from radiation on surface during sunrise - convection breaks up stable nocturnal later and entrains air from above**

## 7. Scatter plot



Daily PM$_{2.5}$ vs visibility based on RH levels during 2013 -2017

RH (%)

Haze definition:
- Visibility <10km & RH <80%.
- Visibility <10km, 80%<RH<95% &
  PM$_{2.5}$ >75 μg/m$^3$.

PM$_{2.5}$ (μg/m$^3$)

Visibility (km)

# 8. Theilsen regression



Visibility at Noibai airport

-284.75 [-296.74, -271.28] units/year ***

# 9. Air-cluster analysis

# 10. CTW longrange/transport

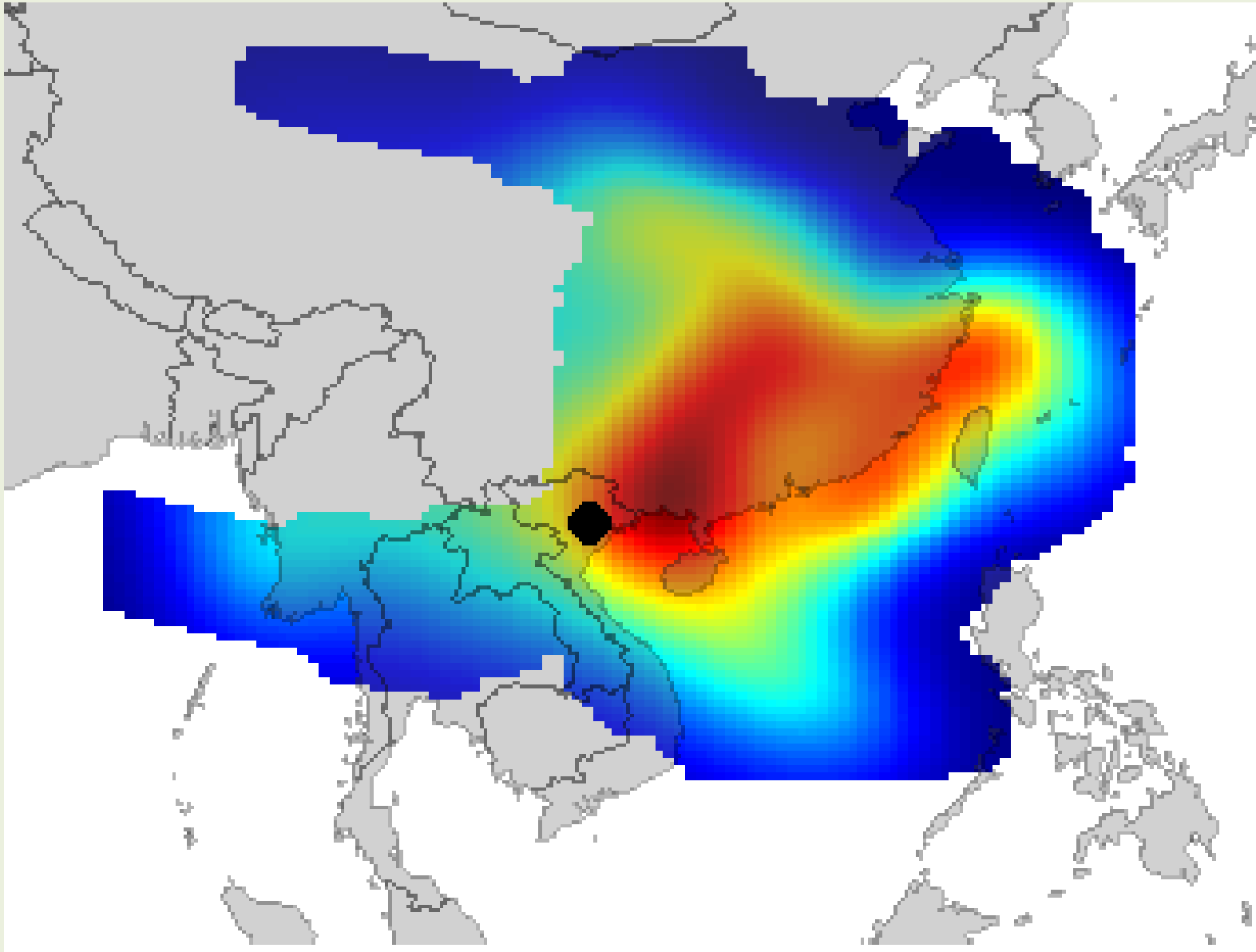# III. Machine learning

**Useful *technique***
1. Factor analysis: PMF, K-mean cluster
2. Decision tree: Random forest, BRT
3. Deep learning: CCN

https://machinelearningcoban.com/
https://rpubs.com/lengockhanhi

Other program: Python

# Weather normalisation using package "*rmweather*"

https://github.com/skgrange/rmweather

Random forest algorithm:
- What is decision tree?
- Random forest is non-linear regression?
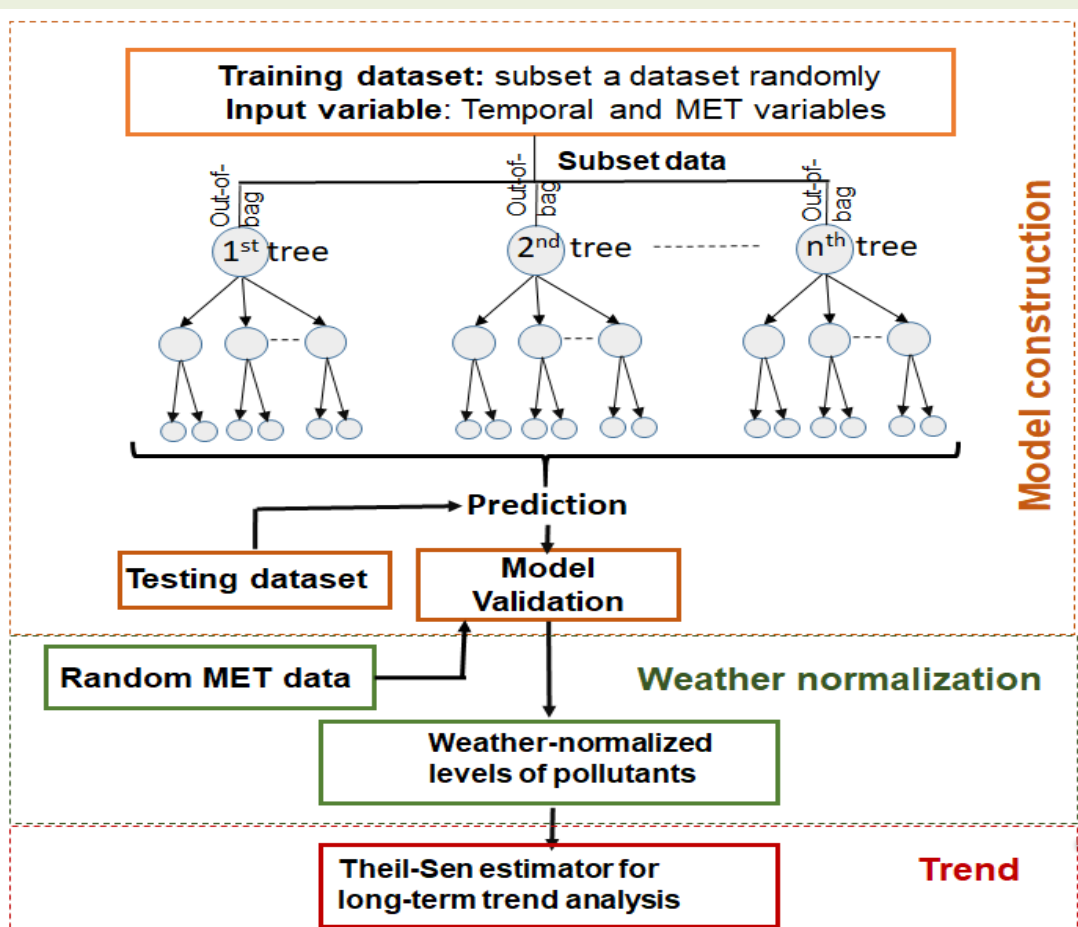- How to select the trees?

# Long-term trend analysis method

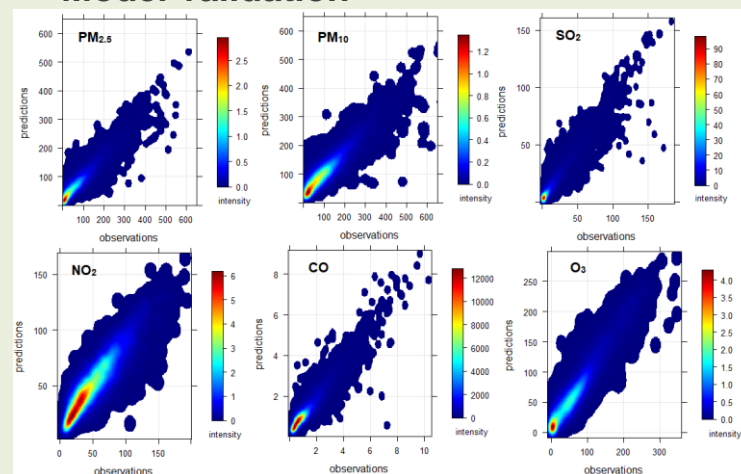The long-term time series of a pollutant can be split into components

$$\ln[C(t)] = C^{LT}(t) + C^{S}(t) + C^{STM}(t) + C^{WH}(t) + C^{WN}(t)$$

*where, LT: long-term component; S: seasonal components; STM: Short-term component; WH: weekend/holiday impact; WN: white-noise is the residual.*
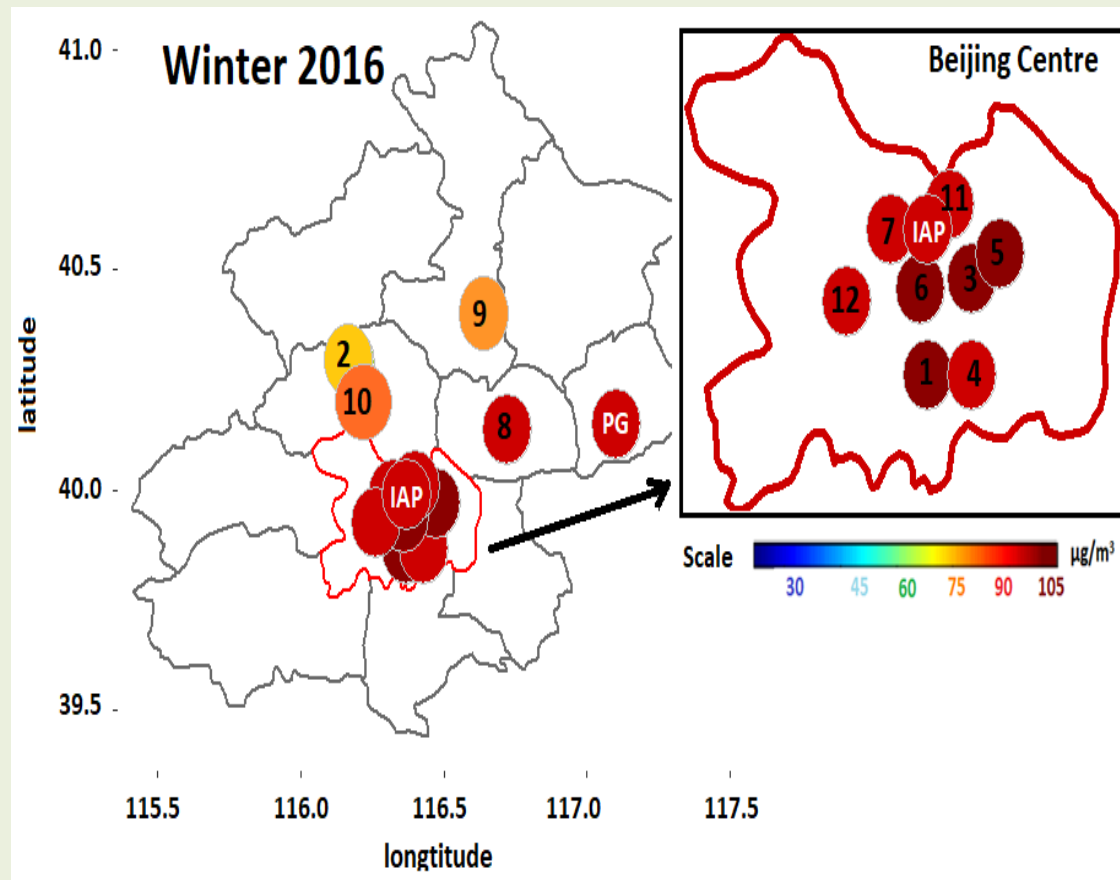
## A decision tree-based random forest technique



**Model validation**



1. **List of policies**
2. **MEIC-emission inventory**

*Refs: Carlaws et al. AE. 2009, Grange et al. ACP 2018, Vu et al . ACPD.2019*
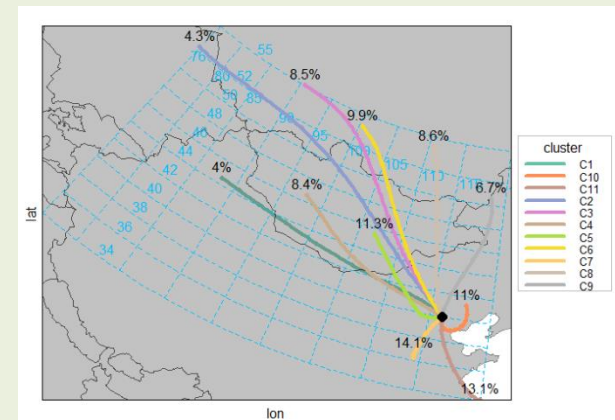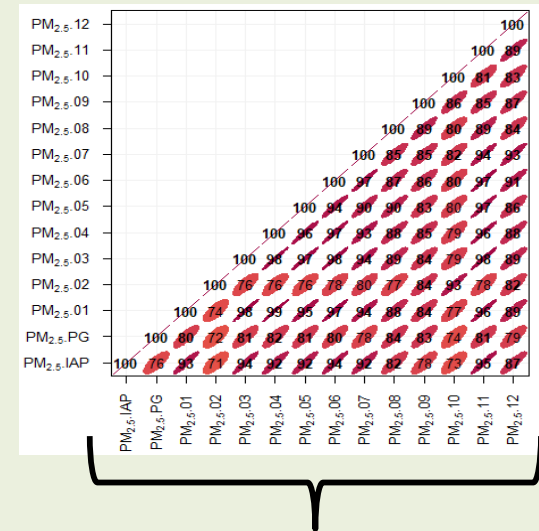
# Input datasets of Air Pollutants in Beijing

Six key pollutants: **PM$_{2.5}$, PM$_{10}$, SO$_2$, NO$_2$, CO, O$_3$**
from 12 national monitoring stations during 2013-2017
& **30-year MET** data sets





**Urban/Suburban/Rural**



**Spatial variation of PM$_{2.5}$ level during APHH winter campaign 2016**

*Refs: Shi et al 2019 ACP*

**30-year MET** data sets
& back trajectories

# Home key messages

- Understand the data first

- Practise basic coding more as you can

- How to use the techniques

## Thank you for your attention