

# Notes on Processing Pollutant Data

This note gives a brief description of how the pollutant data series were constructed from the raw monitor data.

1. Each city has associated with it a number of monitors for a given pollutant. The possible pollutants are PM<sub>10</sub>, PM<sub>2.5</sub>, SO<sub>2</sub>, O<sub>3</sub>, NO<sub>2</sub>, and CO.
2. Let  $X_{j,t}^c$  be the raw pollutant value for monitor  $j$  in city  $c$  at time/day  $t$ . The detrended value  $\tilde{X}_{j,t}^c$  is defined as

$$M_{j,t}^c = \frac{1}{365} \sum_{\ell=-182}^{182} X_{j,t-\ell}^c$$

$$\tilde{X}_{j,t}^c = X_{j,t}^c - M_{j,t}^c$$

The values  $\tilde{X}_{j,t}^c$  are the detrended “residuals” from the raw pollutant series.

3. If a city only has one monitor, then the series  $\tilde{X}_{j,t}^c$  is the final result and is used for analysis.
4. If a city has 2 monitors, then a final  $\bar{X}_t^c$  is computed for each time point  $t$  as

$$\bar{X}_t^c = \frac{1}{2} \left( \tilde{X}_{1,t}^c + \tilde{X}_{2,t}^c \right)$$

If a city has more than 2 monitors, then a 10% trimmed mean of the  $\tilde{X}_{j,t}^c$ ’s is taken for each day. That is, if there are  $J$  monitors in a city, then for each timepoint  $t$

$$\bar{X}_t^c = \text{TrimmedMean}_{10\%} \left[ \tilde{X}_{1,t}^c, \dots, \tilde{X}_{J,t}^c \right]$$

If there are fewer than 10 monitors, the lowest and highest values for each day are still always discarded.

One can see now why the detrending must be done first in Step 2. If a particular monitor has a higher overall level, then it will consistently be discarded when the trimmed mean is taken.

5. The series  $\bar{X}_t^c$  is used as the pollutant measurement for city  $c$  on day  $t$ . In each city dataframe, this series is given the name **\*tmean** where “\*” is either **pm10**, **pm25**, **so2**, **o3**, **no2**, or **co**.
6. The median of the 365-day moving averages are also computed, that is

$$\bar{M}_t^c = \text{Median} \left[ M_{1,t}^c, \dots, M_{J,t}^c \right]$$

This series is given the name **\*mtrend** in each dataframe where “\*” is the name of a pollutant.

The dataframes do not contain the original monitor data, but if one wishes to examine a series that is reminiscent of a true pollutant series, one can add the **\*tmean** series to the **\*mtrend** series. For example, to construct a PM<sub>10</sub> series, one can add the **pm10tmean** and **pm10mtrend** variables.