

Comments on “Filling the gaps of in-situ hourly PM_{2.5} concentration data with the aid of empirical orthogonal function constrained by diurnal cycles”

The method recommended by the authors for the missing value filling to hourly PM_{2.5} data is interesting. It could be useful for relevant study.

Some concerns remain as following, which might be considered to further improve the method.

- (a) Because the PM_{2.5} diurnal variation could vary largely from day to day, is it possible that some typical classification of PM_{2.5} diurnal variation could be established and considered, which should be helpful if one can determine the general pattern of PM_{2.5} diurnal variation for the interested day and make more adequate filling for the missing PM_{2.5} data.
- (b) The PM_{2.5} diurnal variation could be related to some specific meteorological factors as well as their diurnal evolution. Is it possible that the diurnal variation of specific meteorological factors be considered within the authors recommended missing value filling method?
- (c) What is the applicability of the method? Especially for the different spatial distribution of the air quality monitoring stations which are condense over eastern China but sparse over western part of the country.
- (d) In the manuscript, the authors made cross validation for missing value filling for several hours, is it possible that there are missing value for a specific station for one day or several days? If this situation happens, how about the performance of the authors recommended method to make missing value filling?

Some specific comments are also listed below for the authors.

1. Line 60, “data cleaning processes”, consider using more accurate wording to describe what the authors want to mention.
2. Lines 70-71, it is better to directly give the disadvantages of “approaches of ignoring missing values or excluding records on days with missing values”, rather than arbitrarily comment these approaches as “unreasonable”.

3. Table 1, the lines for the references are not quite clear, it is difficult to find which reference is corresponding to which method.
4. Line 152, “ m was defined as the number of stations within 100 km of the target station”, as the authors mentioned about the “significant heterogeneity” of the $PM_{2.5}$ data, is the setting of “100 km” improperly greater in this context? $PM_{2.5}$ concentration can vary largely even within a small area.
Moreover, the air quality monitoring stations are densely distributed over eastern China but sparsely over western part of China. Is there any special consideration should be taken on this issue?
5. The day-to-day $PM_{2.5}$ diurnal variation could vary largely, which depends on whether it is a clean day or a severe polluted day, as well as the various weather conditions. The authors also mentioned this in Lines 302-304. While the method the authors suggested only considers the diurnal variation of one week before and one week after the data missing day to be filled. Is it possible any variety in the diurnal variation of $PM_{2.5}$ can be considered in the recommended method? Also, more detailed classification and establishment of the typical patterns of $PM_{2.5}$ diurnal variation and adequate consideration of this issue could be very helpful to improve the data filling method suggested.
6. Figure 3, it is a little difficult to understand the variables illustrated. The result presented in each panel of the figure seems not match with the caption. The name of the x axis in Figure 3f could be better as “hour”.
7. Figure 4a, the 50th percentile of the mean relative differences generally remains constant around zero, does this mean that the 50th percentile is subjective of less influence from missing values?
8. Figure 6, the reconstructed diurnal $PM_{2.5}$ variation seems to be a smoothed average of the observations near the interested station within a week before and after the interested day, it cannot reconstruct any particular variation of $PM_{2.5}$ such as those at 19:00 local time in Figure 6e and at 08:00-09:00 local time in Figure 6f.
9. Lines 409-411, because of the “significant heterogeneity” of the $PM_{2.5}$ spatial distribution, how about the spatial distribution of the diurnal pattern of $PM_{2.5}$

variation? Is it practical to consider the variability of PM_{2.5} at the stations 100 km away to fill missing value of PM_{2.5}?

10. Do Figure 10a and 10b reflect the same information from different perspectives? Is it possible just keep one figure to discuss the issue?
11. Lines 414-422 and Figure 10, have the authors done data filling for all the available PM_{2.5} data over China with the recommended method? Is the evaluation presented here are based on data filling for the whole dataset of PM_{2.5} available?