

Slower than expected reduction in annual PM_{2.5} in Xi'an revealed by machine learning-based meteorological normalization

Meng Wang¹, Zhuozhi Zhang¹, Qi Yuan¹, Xinwei Li¹, Shuwen Han¹, Yuethang Lam¹, Long Cui², Yu Huang², Junji

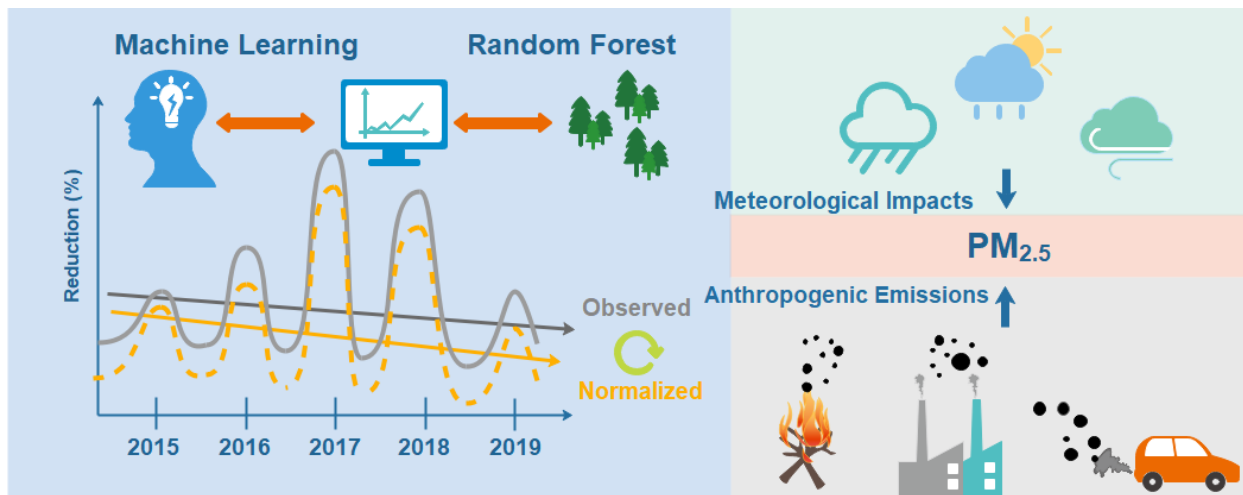
Cao^{2,3}, Shun-cheng Lee^{1, *}

¹Department of Civil and Environmental Engineering, The Hong Kong Polytechnic University, Hung Hom, Hong Kong

²State Key Laboratory of Loess and Quaternary Geology, Institute of Earth Environment, Chinese Academy of Sciences, Xi'an 710061, China

³Key Laboratory of Middle Atmosphere and Global Environment Observation, Institute of Atmospheric Physics, Chinese Academy of Sciences, Beijing 100029, China

Correspondence to: Shun-cheng Lee (shun-cheng.lee@polyu.edu.hk).



Highlights

- Trend analysis of $PM_{2.5}$ over multiple years is complicated due to the impact of meteorology.
- Meteorological normalization was performed using the machine learning algorithm.
- Real trend in $PM_{2.5}$ in a polluted northwest city was revealed after meteorological normalization.
- Reduction rate in the normalized $PM_{2.5}$ over the 5 years was slower than the observed ones.
- Insights into the photochemical and aqueous phase chemistry of secondary $PM_{2.5}$ were gained.

1 **Slower than expected reduction in annual PM_{2.5} in Xi'an**
2 **revealed by machine learning-based meteorological**
3 **normalization**

4 Meng Wang¹, Zhuozhi Zhang¹, Qi Yuan¹, Xinwei Li¹, Shuwen Han¹, Yuethang Lam¹, Long Cui², Yu
5 Huang², Junji Cao^{2,3}, Shun-cheng Lee^{1, *}

6 ¹Department of Civil and Environmental Engineering, The Hong Kong Polytechnic University, Hung
7 Hom, Hong Kong

8 ²State Key Laboratory of Loess and Quaternary Geology, Institute of Earth Environment, Chinese
9 Academy of Sciences, Xi'an 710061, China

10 ³Key Laboratory of Middle Atmosphere and Global Environment Observation, Institute of Atmospheric
11 Physics, Chinese Academy of Sciences, Beijing 100029, China

12

13 *Correspondence to:* Shun-cheng Lee (shun-cheng.lee@polyu.edu.hk).

14

15 **Abstract.** To evaluate the effectiveness of air pollution control policies, trend analysis of the air
16 pollutants is often performed. However, trend analysis of air pollutants over multiple years is complicated
17 by the fact that changes in meteorology over time can also affect the levels of air pollutants in addition
18 to changes in emissions or atmospheric chemistry. To decouple the meteorological effect, this study
19 performed a trend analysis of the hourly fine particulate matter (PM_{2.5}) observed at an urban background
20 site in Xi'an city over 5 years from 2015 to 2019 using the machine learning algorithm. As a novel way
21 of meteorological normalization, the meteorological parameters were used as constant input for 5
22 consecutive years. In this way, the impact of meteorological parameters was excluded, providing insights
23 into the “real” changes in PM_{2.5} due to changes in emission strength or atmospheric chemistry. After
24 meteorological normalization, a decreasing trend of $-3.3\% \text{ year}^{-1}$ ($-1.9 \mu\text{g m}^{-3} \text{ year}^{-1}$) in PM_{2.5} was seen,
25 instead of $-4.4\% \text{ year}^{-1}$ from direct PM_{2.5} observation. Assuming the rate of $-1.9 \mu\text{g m}^{-3} \text{ year}^{-1}$ were kept
26 constant for the next few decades in Xi'an, it would take approximately 25 years (in the year 2045) to
27 reduce the annual PM_{2.5} level to $5 \mu\text{g m}^{-3}$, the new guideline value from World Health Organization. We
28 also show that PM_{2.5} is primarily associated with anthropogenic emissions, which, underwent aqueous
29 phase chemistry in winter and photochemical oxidation in summer as suggested by partial dependence
30 of RH and O_x in different seasons. Therefore, reducing the anthropogenic secondary aerosol precursors
31 at a higher rate, such as NO_x and VOCs is expected to reduce the particulate pollution in this region more
32 effectively than the current $-3.3\% \text{ year}^{-1}$ found in this study.

33

34 **Keywords:** Particulate matter; Secondary aerosol; Theil-Sen estimator; Random forest; Aqueous phase
35 chemistry

36

37 **1 Introduction**

38 Atmospheric particulate matter with a diameter of less than 2.5 μm ($\text{PM}_{2.5}$) is associated with adverse
39 health effects and plays a key role in climate change (Cai et al., 2017; Daellenbach et al., 2020; Wu et
40 al., 2022). Globally, atmospheric $\text{PM}_{2.5}$ is causing millions of premature deaths every year (Burnett et al.,
41 2018; Cohen et al., 2017; Lelieveld et al., 2015). In particular, exposure to high levels of $\text{PM}_{2.5}$ is
42 associated with a high risk of cardiovascular and respiratory disease (Brehmer et al., 2019; Lyu et al.,
43 2018; Yu et al., 2019). The World Health Organization (WHO) recommends an annual $\text{PM}_{2.5}$ level of 10
44 $\mu\text{g m}^{-3}$ (or more recently 5 $\mu\text{g m}^{-3}$) not to be exceeded (WHO, 2006; WHO, 2021). However, it is noted
45 that there is no safe level of $\text{PM}_{2.5}$ below which no adverse health effects would be anticipated (WHO,
46 2006).

47 $\text{PM}_{2.5}$ can be directly emitted from sources of e.g., traffic, industry, and coal combustion; it can also
48 be formed from the oxidation of its precursor gases of e.g., NO_x , SO_2 , volatile organic compounds (VOCs)
49 (Fuzzi et al., 2015; Shrivastava et al., 2017; Zhang et al., 2015), termed as primary and secondary $\text{PM}_{2.5}$,
50 accordingly. To evaluate the effectiveness of air pollution control policies and to further inform policy
51 development, trend analysis of the observed $\text{PM}_{2.5}$ in the ambient environment upon changes in emission
52 and atmospheric chemistry over time is important, especially on a long-term basis e.g., years to decades.
53 However, trend analysis of $\text{PM}_{2.5}$ over multiple years is complicated because changes in meteorology
54 can drive the changes in the observed $\text{PM}_{2.5}$ in addition to changes in emissions or atmospheric chemistry.
55 Therefore, it is essential to decouple the meteorological impact from the observed $\text{PM}_{2.5}$ to see the real
56 changes caused by emission over time with statistical significance.

57 To confirm the changes in pollutant concentration over multiple years with statistical significance, a
58 process called meteorological normalization was proposed by Grange et al. (2018) using the random

59 forest-based machine learning algorithm. The random forest model is computationally efficient and can
60 well predict the PM_{2.5} based on the meteorological parameters (Vu et al., 2019; Zhan et al., 2022; Zhou
61 et al., 2022). By eliminating the effect of meteorological parameters (i.e., after meteorological
62 normalization), insights into the real changes due to emission strength over time can be gained (Grange
63 and Carslaw, 2019; Grange et al., 2021; Grange et al., 2017). Moreover, Qin et al. (2022) show the
64 nonlinear effect of atmospheric variables on the primary and secondary organic aerosol can be well
65 captured by the random forest model. Using the partial dependence algorithm, the emission sources and
66 formation process of PM_{2.5} can be revealed in a complex urban environment (Qin et al., 2022).

67 In China, PM_{2.5} pollution is particularly serious due to rapid economic development, industrialization,
68 and urbanization (Lu et al., 2013; Zhang et al., 2012). To tackle air pollution, many measures have been
69 implemented e.g., the 5-year Clean Air Action Plan and the blue-sky action (Cheng et al., 2019; Wang et
70 al., 2014; Yang et al., 2015). Despite the efforts to reduce the emission, recent studies show the annual
71 PM_{2.5} concentration in Northern China is still far exceeding the WHO guideline values (Chen et al., 2019;
72 Cheng et al., 2019; Vu et al., 2019), highlighting the challenges to improve the air quality in China. As
73 the largest city in northwest China and home to 13 million people as of 2021, Xi'an has suffered severe
74 air pollution over the past decades with PM_{2.5} levels typically higher than in Beijing (Dai et al., 2018;
75 Elser et al., 2016; Huang et al., 2014; Niu et al., 2016). However, compared to Beijing, trend analysis of
76 PM_{2.5} in this highly polluted city is lacking, limiting our understanding of the most recent changes in the
77 evolution of PM_{2.5} over time. In particular, while it is widely acknowledged that 5-year Clean Air Action
78 Plan is contributing to the reduction of PM_{2.5} levels in Beijing (Cheng et al., 2019; Vu et al., 2019), it is
79 unknown if the clean air action is working in Xi'an as significantly.

80 In this study, trend analysis of the hourly PM_{2.5} over 5 years from 2015 to 2019 in Xi'an was performed

81 using the random forest model. The random forest model was used to predict the $PM_{2.5}$ using the
82 meteorological parameters as the model input. Through the comparison of trend analysis before and after
83 meteorological normalization, the effect of meteorological on trend estimates is revealed. Using the
84 partial dependence algorithm, the nonlinear effects of atmospheric variables and gaseous pollutant on
85 $PM_{2.5}$ was evaluated. Finally, implications from trend analysis of $PM_{2.5}$ over the 5 years in Xi'an are
86 discussed.

87 **2 Method**

88 **2.1 Data source**

89 Five years of air quality data (from 2015 to 2019) of the hourly $PM_{2.5}$, NO_2 , SO_2 , O_3 and CO at three
90 national air quality monitoring stations in Xi'an were downloaded from the China National
91 Environmental Monitoring Network website (<https://www.cnemc.cn/>; last access: February 1, 2022). The
92 three sampling sites are all within the urban Xi'an, specifically in three different districts in Xi'an, with
93 GXXQ in Gaoxin District, XZ in Yanta District, and LTQ in Lintong District (Fig. S1). The distance
94 between GXXQ and LTQ sampling site is approximately 40 km, while it is 5 km between GXXQ and
95 XZ (Fig. S1). With such large spatial coverage, the air quality data recorded at the three sampling sites
96 can represent the overall air quality in Xi'an city, one of the most polluted cities in China.

97 Hourly meteorological data including wind speed, wind direction, temperature, relative humidity (RH)
98 recorded at Xi'an Xianyang International Airport were downloaded using the “worldMet” R package
99 (Carslaw, 2017). Planetary boundary layer (PBL) height and atmospheric pressure were obtained from
100 the reanalysis data at 100 m above ground level at the sampling site of GXXQ using the Hybrid Single-
101 Particle Lagrangian Integrated Trajectory (HYSPPLIT) model (Draxler and Rolph, 2003), developed by

102 the National Oceanic and Atmospheric Administration (NOAA). Data were analyzed in RStudio with a
103 series of packages, including “openair”, “normalweather”, and “ggplot2” (Carslaw and Ropkins, 2012;
104 Grange et al., 2018; Vu et al., 2019).

105 **2.2 Random Forest modelling**

106 **2.2.1 Building the Random Forest model**

107 A decision-tree-based random forest model was developed to understand the trend of the observed PM_{2.5}
108 over the 5 years and to gain insights into the formation pathways of PM_{2.5}. Specifically, the random forest
109 model was built to derive the relationship between PM_{2.5} and its predictor features including time
110 variables (date_unix (number of seconds since 1 January 1970), day of the year (day_julian), weekday,
111 and hour of the day), meteorological parameters (wind speed, wind direction, temperature, relative
112 humidity (RH), PBL, and pressure). The time variables act as proxies for emission strength as they vary
113 in time and season.

114 In the RF model, the whole dataset was randomly divided into a training dataset to build the model
115 and a testing dataset to test the model performance. The training dataset was comprised of 80% of the
116 whole dataset, with the testing data (20%) used to validate the models once the forest had been grown.
117 The number of the independent/explanatory variables used to grow a tree was set to three, while the
118 minimum nod-size was set to five, following (Grange et al., 2018). The number of trees within a forest
119 was set to 300. The RF model was built using the latest “rmweather” R package developed by Grange et
120 al. (2018).

121 **2.2.2 Meteorological normalization**

122 PM_{2.5} can be meteorologically normalized by repeatedly (1000 times) re-sampling and predicting using
123 the random models as detailed by Grange et al. (2018). Briefly, PM_{2.5} at a specific measured time point
124 with randomly resampled explanatory variables (except for date_unix) is predicted 1000 times and
125 averaged. For every prediction, the explanatory variables including the time variables (excluding the
126 date_unix variable) and meteorological parameters were randomly selected from the original observation
127 dataset and were subsequently fed to the RF model to predict PM_{2.5} at that particular time point. This is
128 repeated 1000 times, and the 1000 predictions were then averaged, representing “average”
129 meteorological conditions and hence, was regarded as the meteorologically normalized trend. In other
130 words, the meteorological normalized PM_{2.5} (in $\mu\text{g m}^{-3}$) can be thought of as concentrations in “average”
131 or invariant weather conditions. Because the time variables of the hour, weekday, day of the year are also
132 included for normalization, it is not straightforward to investigate the hourly, weekday, seasonal for a
133 comparison with the trend of the observed values.

134 In this study, the meteorological parameters in 2015 were used as the input to predict the PM_{2.5}
135 concentrations in 2016, 2017, 2018, and 2019. In other words, the predicted PM_{2.5} in 2016, 2017, 2018,
136 2019 were the expected PM_{2.5} concentrations under 2015 meteorological conditions. In this way, the
137 predicted PM_{2.5} can be directly compared with the observed PM_{2.5} in terms of hourly, weekday, seasonal
138 variations. Note that only the meteorological parameters (i.e., wind speed, wind direction, temperature,
139 pressure, and PBL) were re-sampled, while the time variables were unchanged. With the predicted PM_{2.5}
140 under the same meteorological conditions from 2015-2019, the behavior of the PM_{2.5} trend due to the
141 changes in emissions or atmospheric chemistry can be revealed.

142 **2.2.3 Partial dependence algorithm**

143 The partial dependence algorithm was applied to assess the nonlinear effect of atmospheric variables,
144 including physical and chemical processes, on the measured PM_{2.5} (Grange and Carslaw, 2019; Grange
145 et al., 2018). The partial dependence algorithm calculated the dependence between the PM_{2.5} and the
146 target atmospheric variables while holding other variables constant at their averages. By targeting all
147 variables one by one, the partial dependence of PM_{2.5} on all considered atmospheric variables was
148 calculated.

149 In this study, the atmospheric variables used as model input included meteorological parameters and
150 gas pollutants. Specifically, the meteorological parameters were RH and temperature which are key
151 atmospheric variables that can influence the physical and chemical processes of PM_{2.5}. For example, high
152 RH may promote aqueous phase chemistry (Duan et al., 2020), while the high temperature may induce
153 high biogenic VOC emissions in summer, key precursor gases for secondary aerosol. Gas pollutants
154 include CO, SO₂, NO₂, O₃, as well as O_x (NO₂ + O₃). CO and SO₂ are indicators of primary emissions,
155 while, O_x is a good surrogate of the oxidizing capability of the atmosphere (Lin et al., 2020). Note that
156 although CO and SO₂ are primary emissions, they are not necessarily local since they can be transported
157 from upwind regions to the receptor sites. The partial dependence algorithm is provided in the
158 “rmweather” package (Grange et al., 2018) in R (version 4.1.2)

159 **2.3 Trend analysis using Theil-Sen estimator**

160 The Theil-Sen regression methodology was applied to investigate the long-term trend of PM_{2.5} before
161 and after the meteorological normalization. The Theil-Sen approach is commonly used for long-term
162 trend analysis and has been detailed in Grange et al. (2018) and Vu et al. (2019). Briefly, the Theil-Sen

163 regression approach accounted for autocorrelation and was used at the 95% confidence level to indicate
164 a significant trend (Grange et al., 2018). The Theil-Sen approach computed the slopes of all possible
165 pairs of $PM_{2.5}$ and took the median values of the slopes, resulting in more conservative confidence
166 intervals for $PM_{2.5}$ trend analysis. The Theil-Sen functions are provided in the “openair” package in R
167 (version 4.1.2) (Carslaw and Ropkins, 2012).

168 **3 Results and Discussion**

169 **3.1 Ambient $PM_{2.5}$ in Xi'an from 2015 to 2019**

170 Figure 1 shows the daily averaged time series of $PM_{2.5}$ over the five years from 2015 to 2019 at the three
171 different sites (i.e., LTQ, XZ, and GXXQ) in Xi'an. The time series of $PM_{2.5}$ at the three sites were very
172 similar with elevated concentrations in winter (spiking over $400 \mu g m^{-3}$) and relatively reduced
173 concentrations in summer ($< 100 \mu g m^{-3}$). Averaged over the five years, $PM_{2.5}$ was 65.1 ± 59.9 (SD) μg
174 m^{-3} at GXXQ, while it was $62.2 \pm 61.2 \mu g m^{-3}$ and $59.3 \pm 58.6 \mu g m^{-3}$ at XZ and LTQ, respectively (Table
175 S1). Despite the large distance between the sampling sites (up to 40 km; Figure S1), the time series of
176 $PM_{2.5}$ at the three sites were highly correlated with correlation coefficient $r > 0.85$ (p -value < 0.01) and
177 slopes close to unity. The good correlation for the observed $PM_{2.5}$ at the three sampling sites suggests the
178 observed $PM_{2.5}$ were due to common pollution sources, simultaneously impacting the air quality over a
179 large area in Xi'an with a diameter of at least 40 km. Due to the similar trend in time series and the
180 slightly high concentration observed at GXXQ, below we focus on the discussion on the air quality data
181 at GXXQ.

182 In terms of annual mean concentration, the $PM_{2.5}$ at GXXQ was $63.6 \mu g m^{-3}$ in 2015. It increased to
183 $74.8 \mu g m^{-3}$ in 2017 then dropped to $58.8 \mu g m^{-3}$ in 2019 (Table S1). Compared to China's national

184 ambient air quality standard (NAAQS-II) of $35 \mu\text{g m}^{-3}$ and the new WHO guideline of $5 \mu\text{g m}^{-3}$ (WHO,
185 2021), the annual mean $\text{PM}_{2.5}$ concentration in Xi'an was approximately substantially (2-7 times) higher,
186 highlighting the poor air quality in this city. Moreover, compared to the trend of $\text{PM}_{2.5}$ in Beijing (Vu et
187 al., 2019), which showed a decreasing trend from $88 \mu\text{g m}^{-3}$ in 2013 to $58 \mu\text{g m}^{-3}$ in 2017, the $\text{PM}_{2.5}$ trend
188 observed in Xi'an is more complicated since the annual $\text{PM}_{2.5}$ concentration increased in 2017 then started
189 decreased afterward. In particular, the number of haze days (defined as daily $\text{PM}_{2.5} > 75 \mu\text{g m}^{-3}$) was 90
190 days (i.e., ~25% of the year or 1 in 4 days; Table S2) in 2015. It increased to 112 days in 2017 then
191 dropped to 86 days in 2019 (Table S2). Most of the haze days occurred in winter, with the average $\text{PM}_{2.5}$
192 concentrations in the range of $67.3\text{-}143 \mu\text{g m}^{-3}$ in winter (Table S3), roughly three times higher than in
193 summer ($24.5\text{-}38.6 \mu\text{g m}^{-3}$).

194 **3.2 Predicted $\text{PM}_{2.5}$ in a good agreement with the observed $\text{PM}_{2.5}$ over 5 years**

195 A decision-tree-based random forest model was trained for the observed $\text{PM}_{2.5}$ with the independent
196 variables including time variables and meteorological parameters as the model input (see the Method
197 section). During the model building, 80% of the dataset was randomly selected as the training dataset,
198 with the rest 20% as the testing dataset. For the training dataset, the predicted $\text{PM}_{2.5}$ was well correlated
199 with the observed $\text{PM}_{2.5}$ with R^2 of 0.99 and slope of 0.93 (Figure 2), while for the testing dataset, the
200 model reproduced the observed $\text{PM}_{2.5}$ reasonably well with R^2 of 0.93 and slope of 0.84. The slope of
201 0.84-0.93 for the testing dataset suggested the model tended to underestimate the $\text{PM}_{2.5}$ by 7-16%.
202 Nevertheless, the high R^2 values (0.93-0.99) for both the training and testing dataset suggest the random
203 forest grown in this study had a strong explanatory ability for $\text{PM}_{2.5}$.

204 The good performance of the random forest model was partly due to the strong seasonality of the $\text{PM}_{2.5}$

205 which was well captured by the model. Specifically, the time variable (i.e., day of the year or day_Julian
206 (1-365)) was the most important variable for PM_{2.5} explanation in the random forest model (Figure S2).
207 Partial dependence on the time variable of day_julian shows the elevated PM_{2.5} concentrations (> 75 µg
208 m⁻³) were associated with day 1-50 and day 300-365 in the year (Figure S3), consistent with the fact that
209 haze pollution occurred mostly in winter. In contrast, the time variable of weekday and hour were of less
210 importance. Partial dependence plots on the weekday do not present a clear weekday/weekends
211 difference (Figure S3). Given that traffic is usually heavier during weekdays than weekends, the lack of
212 weekday and weekends pattern suggests traffic was not the major source of PM_{2.5}. Consistently, the
213 partial dependence on the time variable of hour shows no rush hour peaks (Figure S3). Instead, elevated
214 PM_{2.5} concentrations were found to be occurring mostly in the night till the next morning (Figure S3).

215 Among the meteorological parameters, temperature was the most important parameter, followed by
216 pressure, relative humidity, boundary layer height, wind speed and wind direction (Figure S2). In
217 particular, the low temperature (< 10 °C), low pressure (< 880 hpa), high relative humidity (50-90%),
218 low boundary layer height (< 500 m), low wind speed (< 3 m s⁻¹) in north-easterly wind were associated
219 with high PM_{2.5} concentrations (Figure S3). These meteorological parameters created a stable atmosphere
220 with poor dispersion conditions, causing the build-up of PM_{2.5}. Note that at higher relative humidity (95-
221 100%), precipitation events likely caused the wet deposition of PM_{2.5}, leading to lower concentrations as
222 a result (Figure S3). Therefore, the random forest model captured the general predictions and processes
223 that were associated with the ambient PM_{2.5} concentrations, confirming the strong explanatory power of
224 the model.

225 3.3 Trend analysis before and after meteorological normalization

226 [Figure 3](#) shows the yearly averaged $PM_{2.5}$ concentration before and after meteorological normalization.

227 It shows that $PM_{2.5}$ concentration in 2019 would have been higher if under the 2015 meteorological

228 conditions, while $PM_{2.5}$ concentrations in 2017 would have been lower if under the same 2015

229 meteorological conditions. Specifically, the observed $PM_{2.5}$ concentrations (i.e., before meteorological

230 normalization) were $63.6 \mu\text{g m}^{-3}$, $66.5 \mu\text{g m}^{-3}$, $74.8 \mu\text{g m}^{-3}$, $61.2 \mu\text{g m}^{-3}$, $58.8 \mu\text{g m}^{-3}$, in 2015, 2016, 2017,

231 2018, 2019, respectively. After meteorological normalization, the predicted $PM_{2.5}$ concentrations from

232 2015 to 2019 were $63.4 \mu\text{g m}^{-3}$, $66.8 \mu\text{g m}^{-3}$, $72.3 \mu\text{g m}^{-3}$, $64.1 \mu\text{g m}^{-3}$, $60.9 \mu\text{g m}^{-3}$, respectively. This can

233 be translated to percentages differences of -0.3% , 0.4% , -3.3% , 4.7% , 3.6% by comparing the predicted

234 and observed $PM_{2.5}$. The percentage differences may appear small (from -3.3% to 4.7%) when compared

235 with the observed $PM_{2.5}$. However, in terms of the yearly $PM_{2.5}$ trend analysis which is usually on the

236 scale of 1-10% ([Vu et al., 2019](#)), the changes in $PM_{2.5}$ due to different meteorological conditions may

237 have a big impact. Below, we discuss the effect of meteorological normalization on $PM_{2.5}$ trend analysis.

238 [Figure 4](#) shows the trend analysis of monthly averaged $PM_{2.5}$ before and after meteorological

239 normalization using the same Theil-Sen algorithm (see Sect. 2.4). The temporal variations of the monthly

240 average $PM_{2.5}$ for both cases do not show a smooth trend from 2016 to 2019 because of the spikes during

241 pollution events in cold seasons, consistent with the daily average $PM_{2.5}$ as shown in [Figure 1](#). Using the

242 Theil-Sen estimator, the observed $PM_{2.5}$ (before meteorological normalization) shows a trend of -4.4%

243 year^{-1} or $-2.6 \mu\text{g m}^{-3}$, while it shows a less negative trend of $-3.3\% \text{ year}^{-1}$ ($-1.9 \mu\text{g m}^{-3}$) after

244 meteorological normalization. However, both show a large range in terms of 95% confidence level, from

245 $-10.76\% \text{ year}^{-1}$ to $6.1\% \text{ year}^{-1}$ for the observed $PM_{2.5}$ and slightly positive values from $-9.41\% \text{ year}^{-1}$ to

246 $6.7\% \text{ year}^{-1}$ for the meteorological normalized $PM_{2.5}$. The large range of the 95% confidence level was

247 due to the fact that the pollution events in cold seasons do not appear abating in terms of the magnitude
248 of the PM_{2.5} levels and the duration of the pollution. Nevertheless, compared to the observed PM_{2.5} trend,
249 the slightly positive trend of PM_{2.5} after meteorological normalization results suggest that the effect of
250 emission reduction was contributing less to the improvement of air quality in Xi'an. This is in great
251 contrast to the findings in Beijing, where emission reductions were found to cause a larger reduction in
252 PM_{2.5} after meteorological normalization although with different normalization methodology over
253 different years (2013-2017; (Vu et al., 2019)).

254 **3.4 Formation process of haze pollution**

255 As discussed above, the most severe pollution events occurred in the winter months (December, January,
256 and February) over the 5 years from 2016 to 2019. To gain insights into the formation processes of haze
257 pollution in winter, gas pollutants of CO, NO₂, O₃, SO₂, and O_x were fed into the random forest in addition
258 to the meteorological parameters of RH, temperature, wind speed, and wind direction. As a comparison,
259 a similar analysis was also performed during the summer months (June, July, and August). Note that
260 because we focused on only one season, time variables were not considered. Additionally, a multi-linear
261 regression model was also performed between PM_{2.5} and these gas pollutants. The results show the
262 correlation determination R² for the random forest model (0.64-0.71) was significantly higher than for
263 the multi-linear regression model (< 0.4; Table S4). Therefore, the random forest model can provide
264 higher accuracies of the PM_{2.5} prediction than the multi-linear regression model.

265 [Figure 5](#) shows the importance parameters during the model training process for the winter and
266 summer PM_{2.5}. In both winter and summer, CO is the most important variable in explaining the observed
267 PM_{2.5} ([Figure 5](#)). Given that CO is the by-product of incomplete combustion, the strong importance of

268 CO in explaining $PM_{2.5}$ suggests $PM_{2.5}$ was primarily associated with anthropogenic combustion sources
269 including both direct emission and/or secondary formation from anthropogenic precursor gases. However,
270 as discussed above, the good time series correlation between the three sampling sites (despite a distance
271 of 40 km) suggests the observed $PM_{2.5}$ were regionally relevant rather than local pollution events, and,
272 therefore, the observed $PM_{2.5}$ was likely associated with the secondary formation during transport. Indeed,
273 recent studies of $PM_{2.5}$ source apportionment highlight secondary aerosols are the major component of
274 $PM_{2.5}$ instead of primary emission in Xi'an (Duan et al., 2020; Elser et al., 2016; Zhong et al., 2020).

275 [Figure 5](#) also shows, while RH was the second most important parameter in winter, O_x ($NO_2 + O_3$)
276 was the second most important variable in summer. The partial dependence plot shows that high RH was
277 associated with high $PM_{2.5}$ in winter ([Figure 6](#)), while high O_x was associated with high $PM_{2.5}$ in summer.
278 Recent studies show RH can promote the aqueous formation of secondary aerosol including sulfate and
279 oxygenated organic aerosol, while these secondary aerosols together often contribute over half of the
280 $PM_{2.5}$ mass (Elser et al., 2016; Zhong et al., 2020). In this study, the aqueous phase chemistry is reflected
281 by the high importance of RH in winter, with the overall $PM_{2.5}$ showing a positive response to RH in
282 winter. In contrast, O_x is a good indicator of photochemical chemistry, which is more important in
283 summer than in winter as reflected by the importance of O_x in summer ([Figure 5](#)), showing a positive
284 response to O_x ([Figure 6](#)). As a comparison, O_x in winter was the least important gaseous variable in
285 winter, suggesting photochemical chemistry was less significant than RH-promoted aqueous phase
286 chemistry. In summer, RH was still the fourth important variable after O_3 , implying aqueous phase
287 chemistry could also be the major pathway for secondary aerosol formation. Consistently, Duan et al.
288 (2020) shows a large formation of secondary aerosol formation during fog-rain days in summer Xi'an.

289 **4 Discussion**

290 Using the random forest model, we show that the 5-year hourly $PM_{2.5}$ measured at the suburban site in
291 Xi'an from 2015-2019 was reproduced well by feeding the meteorological parameters and time variables
292 into the model. Meteorological parameters can affect the dispersion conditions and/or atmospheric
293 chemistry of the ambient $PM_{2.5}$, while the time variables act as proxies for emission strength as they vary
294 in terms of hour, day, season, and year. Assuming the meteorological parameters were the same
295 throughout the 5 years (i.e., normalization), we can exclude the impact of meteorological parameters,
296 providing insights into the “real” changes in $PM_{2.5}$ due to changes in emission strength or atmospheric
297 chemistry. After meteorological normalization, we show that the $PM_{2.5}$ concentration in 2019 would have
298 been higher, while $PM_{2.5}$ concentrations in 2017 would have been lower if under the same meteorological
299 conditions as in 2015. As a result, a decreasing trend of -3.3% year⁻¹ in $PM_{2.5}$ after meteorological
300 normalization was seen, instead of -4.4% from direct $PM_{2.5}$ observation. The “real” decreasing rate of
301 -3.3% year⁻¹ for $PM_{2.5}$ in Xi'an was roughly half of the values (-7.8% year⁻¹) reported in Beijing over
302 the year of 2013-2017 (Vu et al., 2019). Assuming the rate of -3.3% year⁻¹ or $1.9\ \mu\text{g m}^{-3}$ year⁻¹ were kept
303 constant for the next few decades in Xi'an, it would take approximately 25 years (in the year 2045) to
304 reduce the yearly $PM_{2.5}$ concentration to $10\ \mu\text{g m}^{-3}$, the guideline value from WHO. Therefore, more
305 efforts need to be taken to reduce the $PM_{2.5}$ pollution in this inland city, which is the large northwestern
306 city in China, home to over 10 million in northwest China.

307 We also show that the non-linear effect of atmospheric variables on $PM_{2.5}$ can be captured by the
308 random forest model as opposed to the multi-linear regression. Different from the multi-linear regression
309 model, the random forest model also provides insights into the relative importance of the atmospheric
310 variable. In particular, we show that in both winter and summer, CO is the most important variable,

311 suggesting the observed $PM_{2.5}$ is primarily associated with anthropogenic emissions, which, undergoes
312 aqueous phase chemistry in winter and photochemical oxidation in summer as suggested by importance
313 of RH and O_x , the second most important variable, accordingly, after CO. Given that the time series of
314 $PM_{2.5}$ are well correlated at the three sampling sites, despite a distance of 40 km apart, the secondary
315 formation pathways, which is different in different seasons, play an important role covering a large area
316 in Xi'an. As a result of secondary formation, the difference in $PM_{2.5}$ concentration at the three sampling
317 sites is marginal. Therefore, reducing the anthropogenic secondary aerosol precursors at a higher rate,
318 such as NO_x and VOCs is expected to reduce the particulate pollution in this region at a faster pace than
319 the current -3.3% year⁻¹ found in this study.

320 **5 Conclusion**

321 In this study, trend analysis of the hourly fine particulate matter ($PM_{2.5}$) observed at an urban background
322 site in Xi'an city over 5 years from 2015 to 2019 was performed using the machine learning algorithm -
323 random forest model. To decouple the meteorological effect, the meteorological parameters were
324 assumed the same throughout the 5 years. In this way, the impact of meteorological parameters was
325 excluded, providing insights into the "real" changes in $PM_{2.5}$ due to changes in emission strength or
326 atmospheric chemistry over 5 years. After meteorological normalization, the "real" decreasing trend of
327 -3.3% year⁻¹ in $PM_{2.5}$ after meteorological normalization was roughly 30% higher than the trend of -4.4%
328 year⁻¹ from direct $PM_{2.5}$ observation. Therefore, meteorological normalization made the decreasing trend
329 of $PM_{2.5}$ less significant. The "real" decreasing rate of -3.3% year⁻¹ for $PM_{2.5}$ in Xi'an was roughly half
330 of the values (-7.8% year⁻¹) reported in Beijing over the year of 2013-2017, suggesting the air quality
331 control measures were less effective in this region. To take the decreasing trend into context, we assumed

332 the rate of $-3.3\% \text{ year}^{-1}$ or $1.9 \mu\text{g m}^{-3} \text{ year}^{-1}$ were kept constant for the next few decades in Xi'an. Then,
333 it would take 25 years (in the year 2045) to reduce the yearly $\text{PM}_{2.5}$ concentration to $10 \mu\text{g m}^{-3}$. Through
334 relative importance analysis and partial dependence algorithm, the observed $\text{PM}_{2.5}$ was found to be
335 primarily associated with anthropogenic emissions, which, underwent aqueous phase chemistry in winter
336 and photochemical oxidation in summer. Therefore, reducing the anthropogenic secondary aerosol
337 precursors at a higher rate, such as NO_x and VOCs is expected to reduce the particulate pollution in this
338 region more efficiently. This study provides a robust trend analysis in $\text{PM}_{2.5}$ over 5 years in a highly
339 polluted but less studied city in northwest China, providing high certainty that the real trend is less
340 significant under the current control measures than observed, requiring stricter policies controlling the
341 emission of precursor gases from anthropogenic activities.

342

343

344 **Associate content**

345 Supporting Information

346 Supplementary figures (Fig. S1-S2) and Table S1-S4.

347 **Credit authorship contribution statement**

348 Meng Wang: designed the study, conducted data analysis, prepared the manuscript with contributions
349 from all co-authors.

350 Zhuozhi Zhang: Formal analysis, Writing, Review and Editing.

351 Qi Yuan: Formal analysis, Methodology.

352 Xinwei Li: Investigation, Methodology.
353 Shuwen Han: Validation, Investigation.
354 Yuethang Lam: Formal analysis.
355 Long Cui: Formal analysis, Investigation.
356 Yu Huang: Writing, Review and Editing.
357 Shun-cheng Lee: Writing - review and editing, Funding acquisition, Supervision.

358 **Declaration of competing interest**

359 The authors declare that they have no conflicting interests.

360 **Acknowledgements**

361 This work was supported by the Environment and Conservation Fund - Environmental Research,
362 Technology Demonstration and Conference Projects (ECF 63/2019), the RGC Theme-based Research
363 Scheme (T24-504/17-N), the RGC Theme-based Research Scheme (T31-603/21-N).

364 **References**

365 Brehmer C, Lai A, Clark S, Shan M, Ni K, Ezzati M, et al. The Oxidative Potential of Personal and
366 Household PM_{2.5} in a Rural Setting in Southwestern China. *Environmental Science & Technology*
367 2019; 53: 2788-2798.
368 Burnett R, Chen H, Szyszkowicz M, Fann N, Hubbell B, Pope CA, et al. Global estimates of mortality
369 associated with long-term exposure to outdoor fine particulate matter. *Proceedings of the National*
370 *Academy of Sciences* 2018; 115: 9592-9597.

371 Cai W, Li K, Liao H, Wang H, Wu L. Weather conditions conducive to Beijing severe haze more frequent
372 under climate change. *Nature Climate Change* 2017; 7: 257.

373 Carslaw DC. Worldmet: Import Surface Meteorological Data from NOAA Integrated Surface Database
374 (ISD), available at: <http://github.com/davidcarslaw/> (last access: Feb 01, 2022), 2017.

375 Carslaw DC, Ropkins K. openair — An R package for air quality data analysis. *Environmental Modelling*
376 & *Software* 2012; 27-28: 52-61.

377 Chen D, Liu Z, Ban J, Zhao P, Chen M. Retrospective analysis of 2015–2017 wintertime PM_{2.5} in China:
378 response to emission regulations and the role of meteorology. *Atmospheric Chemistry and Physics*
379 2019; 19: 7409-7427.

380 Cheng J, Su J, Cui T, Li X, Dong X, Sun F, et al. Dominant role of emission reduction in PM_{2.5} air quality
381 improvement in Beijing during 2013–2017: a model-based decomposition analysis. *Atmospheric*
382 *Chemistry and Physics* 2019; 19: 6125-6146.

383 Cohen AJ, Brauer M, Burnett R, Anderson HR, Frostad J, Estep K, et al. Estimates and 25-year trends of
384 the global burden of disease attributable to ambient air pollution: an analysis of data from the Global
385 Burden of Diseases Study 2015. *The Lancet* 2017; 389: 1907-1918.

386 Daellenbach KR, Uzu G, Jiang J, Cassagnes L-E, Leni Z, Vlachou A, et al. Sources of particulate-matter
387 air pollution and its oxidative potential in Europe. *Nature* 2020; 587: 414-419.

388 Dai Q, Bi X-H, Liu B, Li L, Ding J, Song W, et al. Chemical nature of PM_{2.5} and PM₁₀ in Xi'an, China:
389 Insights into primary emissions and secondary particle formation. *Environmental pollution (Barking,*
390 *Essex : 1987)* 2018; 240: 155-166.

391 Draxler RR, Rolph GD. HYSPLIT (HYbrid Single-Particle Lagrangian Integrated Trajectory) model
392 access via NOAA ARL READY website (www.arl.noaa.gov/ready/hysplit4.html). NOAA Air

393 Resources Laboratory, Silver Spring. Md, 2003.

394 Duan J, Huang R-J, Gu Y, Lin C, Zhong H, Wang Y, et al. The formation and evolution of secondary
395 organic aerosol during summer in Xi'an: Aqueous phase processing in fog-rain days. *Science of The
396 Total Environment* 2020; 144077.

397 Elser M, Huang RJ, Wolf R, Slowik JG, Wang Q, Canonaco F, et al. New insights into PM_{2.5} chemical
398 composition and sources in two major cities in China during extreme haze events using aerosol mass
399 spectrometry. *Atmospheric Chemistry and Physics* 2016; 16: 3207-3225.

400 Fuzzi S, Baltensperger U, Carslaw K, Decesari S, Denier Van Der Gon H, Facchini M, et al. Particulate
401 matter, air quality and climate: lessons learned and future needs. *Atmospheric Chemistry and
402 Physics* 2015; 15: 8217-8299.

403 Grange SK, Carslaw DC. Using meteorological normalisation to detect interventions in air quality time
404 series. *Science of The Total Environment* 2019; 653: 578-588.

405 Grange SK, Carslaw DC, Lewis AC, Boleti E, Hueglin C. Random forest meteorological normalisation
406 models for Swiss PM₁₀ trend analysis. *Atmospheric Chemistry and Physics* 2018; 18: 6223-6239.

407 Grange SK, Lee JD, Drysdale WS, Lewis AC, Hueglin C, Emmenegger L, et al. COVID-19 lockdowns
408 highlight a risk of increasing ozone pollution in European urban areas. *Atmospheric Chemistry and
409 Physics* 2021; 21: 4169-4185.

410 Grange SK, Lewis AC, Moller SJ, Carslaw DC. Lower vehicular primary emissions of NO₂ in Europe
411 than assumed in policy projections. *Nature Geoscience* 2017; 10: 914-918.

412 Huang R-J, Zhang Y, Bozzetti C, Ho K-F, Cao J-J, Han Y, et al. High secondary aerosol contribution to
413 particulate pollution during haze events in China. *Nature* 2014; 514: 218.

414 Lelieveld J, Evans JS, Fnais M, Giannadaki D, Pozzer A. The contribution of outdoor air pollution

415 sources to premature mortality on a global scale. *Nature* 2015; 525: 367-371.

416 Lin C, Huang R-J, Xu W, Duan J, Zheng Y, Chen Q, et al. Comprehensive Source Apportionment of
417 Submicron Aerosol in Shijiazhuang, China: Secondary Aerosol Formation and Holiday Effects.
418 *ACS Earth and Space Chemistry* 2020; 4: 947-957.

419 Lu Q, Zheng J, Ye S, Shen X, Yuan Z, Yin S. Emission trends and source characteristics of SO₂, NO_x,
420 PM₁₀ and VOCs in the Pearl River Delta region from 2000 to 2009. *Atmospheric Environment* 2013;
421 76: 11-20.

422 Lyu Y, Guo H, Cheng T, Li X. Particle Size Distributions of Oxidative Potential of Lung-Deposited
423 Particles: Assessing Contributions from Quinones and Water-Soluble Metals. *Environmental*
424 *Science & Technology* 2018; 52: 6592-6600.

425 Niu X, Cao J, Shen Z, Ho SSH, Tie X, Zhao S, et al. PM_{2.5} from the Guanzhong Plain: Chemical
426 composition and implications for emission reductions. *Atmospheric Environment* 2016; 147: 458-
427 469.

428 Qin Y, Ye J, Ohno P, Liu P, Wang J, Fu P, et al. Assessing the Nonlinear Effect of Atmospheric Variables
429 on Primary and Oxygenated Organic Aerosol Concentration Using Machine Learning. *ACS Earth*
430 *and Space Chemistry* 2022; 6: 1059-1066.

431 Shrivastava M, Cappa CD, Fan J, Goldstein AH, Guenther AB, Jimenez JL, et al. Recent advances in
432 understanding secondary organic aerosol: Implications for global climate forcing. *Reviews of*
433 *Geophysics* 2017; 55: 509-559.

434 Vu TV, Shi Z, Cheng J, Zhang Q, He K, Wang S, et al. Assessing the impact of clean air action on air
435 quality trends in Beijing using a machine learning technique. *Atmospheric Chemistry and Physics*
436 2019; 19: 11303-11314.

437 Wang S, Xing J, Zhao B, Jang C, Hao J. Effectiveness of national air pollution control policies on the air
438 quality in metropolitan areas of China. *Journal of Environmental Sciences* 2014; 26: 13-22.

439 WHO. Air quality guidelines: global update 2005: World Health Organization, 2006.

440 WHO. https://www.who.int/health-topics/air-pollution#tab=tab_1 (last access: 01 March 2022), 2021.

441 Wu D, Zheng H, Li Q, Jin L, Lyu R, Ding X, et al. Toxic potency-adjusted control of air pollution for
442 solid fuel combustion. *Nature Energy* 2022; 7: 194-202.

443 Yang Z, Wang H, Shao Z, Muncrief R. Review of Beijing's Comprehensive motor vehicle emission
444 Control program, White Paper, available at:
445 https://theicct.org/sites/default/files/publications/Beijing_Emission_Control_Programs_201511.pdf
446 f (last access: 01 March 2022), 2015.

447 Yu S, Liu W, Xu Y, Yi K, Zhou M, Tao S, et al. Characteristics and oxidative potential of atmospheric
448 PM_{2.5} in Beijing: Source apportionment and seasonal variation. *Science of The Total Environment*
449 2019; 650: 277-287.

450 Zhan J, Liu Y, Ma W, Zhang X, Wang X, Bi F, et al. Ozone formation sensitivity study using machine
451 learning coupled with the reactivity of volatile organic compound species. *Atmospheric*
452 *Measurement and Technology* 2022; 15: 1511-1520.

453 Zhang Q, He K, Huo H. Policy: cleaning China's air. *Nature* 2012; 484: 161.

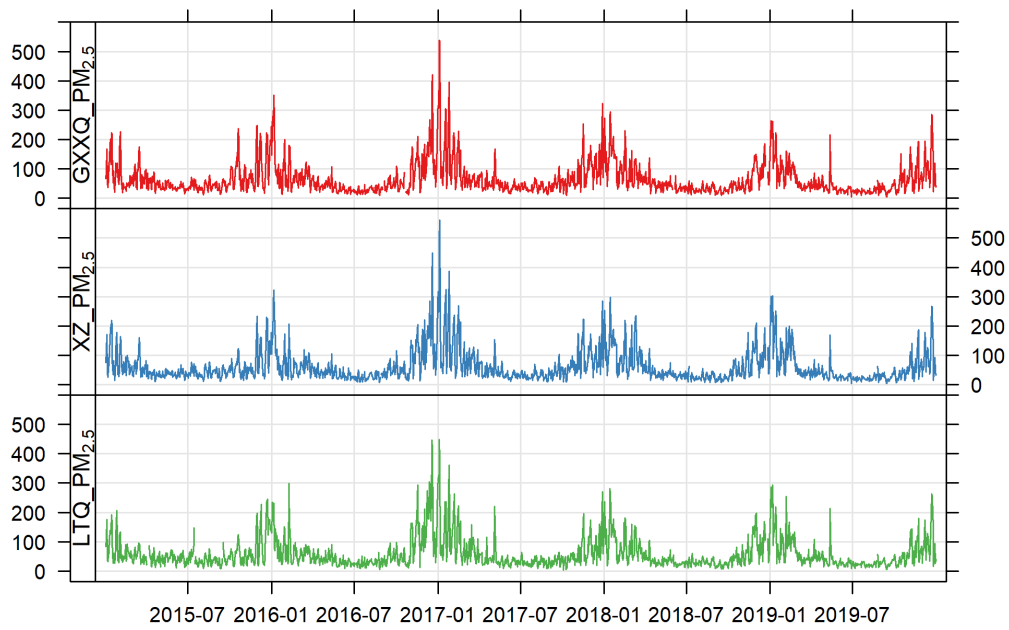
454 Zhang R, Wang G, Guo S, Zamora ML, Ying Q, Lin Y, et al. Formation of Urban Fine Particulate Matter.
455 *Chemical Reviews* 2015; 115: 3803-3855.

456 Zhong H, Huang R-J, Duan J, Lin C, Gu Y, Wang Y, et al. Seasonal variations in the sources of organic
457 aerosol in Xi'an, Northwest China: The importance of biomass burning and secondary formation.
458 *Science of The Total Environment* 2020: 139666.

459 Zhou W, Lei L, Du A, Zhang Z, Li Y, Yang Y, et al. Unexpected increases of severe haze pollution during
460 the post COVID-19 period: effects of emissions, meteorology, and secondary production. Journal
461 of Geophysical Research: Atmospheres 2022: e2021JD035710.

462

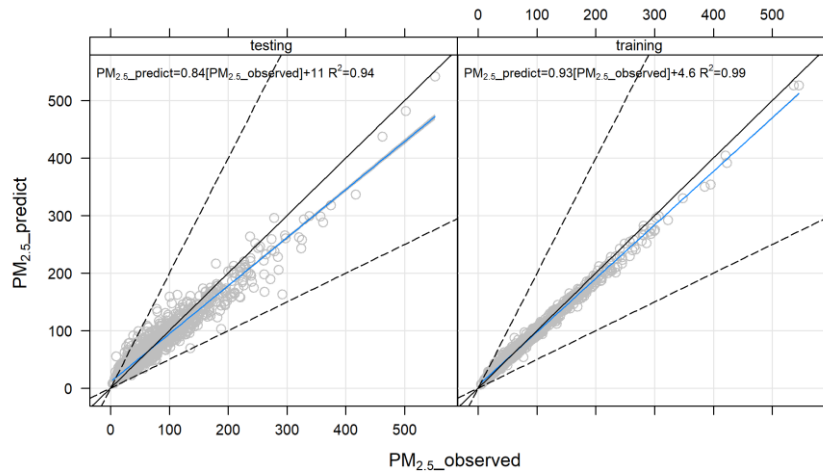
463



464

465 **Figure 1** Time series of the daily averaged PM_{2.5} (in $\mu\text{g m}^{-3}$) at the three sampling sites of GXXQ, XZ,
466 and LTQ, with a distance of up to 40 km apart. A map of the three sampling sites is provided in Fig. S1.

467

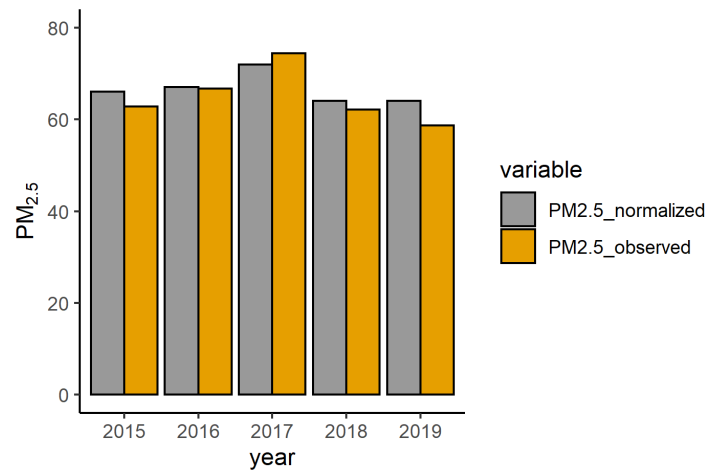


468

469 **Figure 2** Scatter plots of predicted and measured PM_{2.5} concentrations (in µg m⁻³) in the train set and test

470 set. Also shown are the linear correlation and R².

471



472

473 **Figure 3** Annual PM_{2.5} concentration (in µg m⁻³) before (i.e., the observed PM_{2.5}) and after

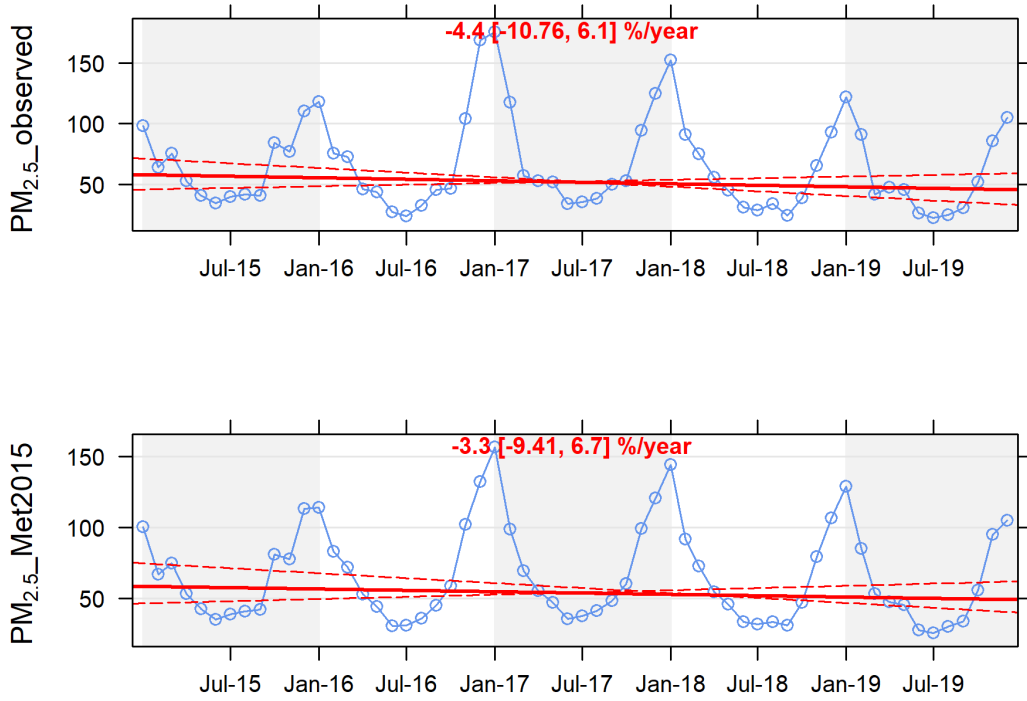
474 meteorological normalization (i.e., the normalized PM_{2.5}).

475

476

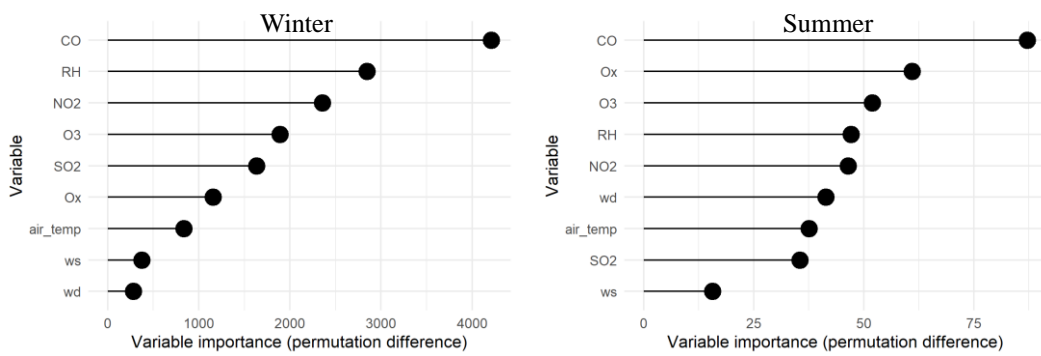
477

478



479

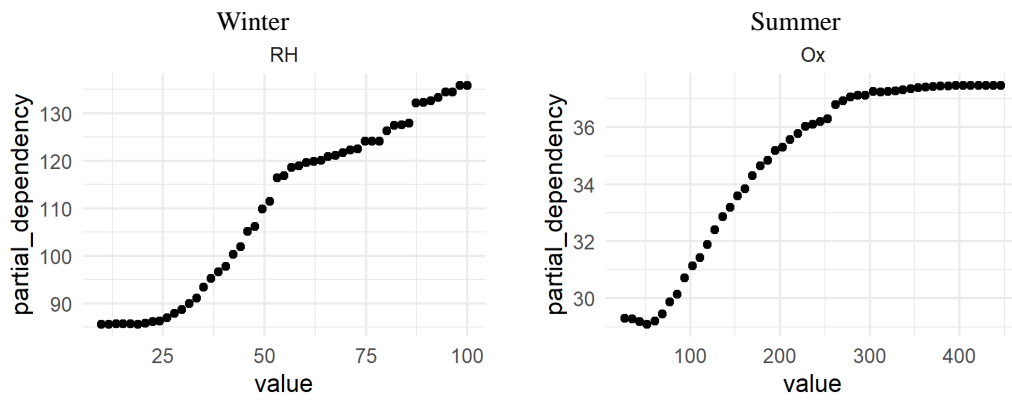
480 **Figure 4** Monthly averaged PM_{2.5} before (top panel) and after (bottom panel) meteorological
481 normalization. The red line represents the trend analysis of PM_{2.5} using the Theil-Sen estimator, with the
482 dotted red line representing the 95% confidence level.



483

484 **Figure 5** Variable importance for the random forest model built for the hourly PM_{2.5} during winter (left
485 panel) and summer (right panel) from 2015 to 2019. The variables include gas pollutants and
486 meteorological parameters.

487



488

489 **Figure 6** Partial dependence plots of the RH in winter (left panel) and O_x in summer (right panel) for the
 490 random forest model built for $PM_{2.5}$.

491



Click here to access/download
Supplementary Material
Supplementary.docx



Credit authorship contribution statement

Meng Wang: designed the study, conducted data analysis, prepared the manuscript with contributions from all co-authors.

Zhuozhi Zhang: Formal analysis, Writing, Review and Editing.

Qi Yuan: Formal analysis, Methodology.

Xinwei Li: Investigation, Methodology.

Shuwen Han: Validation, Investigation.

Yuethang Lam: Formal analysis.

Long Cui: Formal analysis, Investigation.

Yu Huang: Writing, Review and Editing.

Shun-cheng Lee: Writing - review and editing, Funding acquisition, Supervision.