

1 Highlights

2 **Relative Importance of Meteorological Variables on Air Quality and Role of Boundary Layer** 3 **Height**

4 Yaxuan Huang, Bin Guo, Haoxuan Sun, Huijie Liu, Song Xi Chen

- 5 • Stable orders of meteorological effects on six air pollutants in North China.
- 6 • Boundary layer height's effects can be modeled by the surface meteorological variables.
- 7 • $PM_{2.5}$, PM_{10} , SO_2 and CO shared a common order of variable importance.

Relative Importance of Meteorological Variables on Air Quality and Role of Boundary Layer Height

Yaxuan Huang^a, Bin Guo^{b,*}, Haoxuan Sun^c, Huijie Liu^d and Song Xi Chen^{e,**}

^aYuanpei College, Peking University, Beijing, 100871, China

^aDepartment of Statistics, University of California, Berkeley, CA, 94720, USA

^bCenter of Statistical Research and School of Statistics, Southwestern University of Finance and Economics, Chengdu, 611130, China

^cCenter for Data Science, Peking University, Beijing, 100871, China

^dThe Experimental High School Attached To Beijing Normal University, Beijing, 100032, China

^eSchool of Mathematical Sciences, Guanghua School of Management and Center for Statistical Science, Peking University, Beijing, 100871, China

ARTICLE INFO

Keywords:

Air Quality Assessment

Statistical Learning

Meteorological Effects

Variable Selection

ABSTRACT

To gain insight on the meteorological effects of air pollution, we study the relative importance of surface meteorological variables and boundary layer height (BLH) on six major air pollutants according to the orders of variables being selected in a forward variable selection algorithm. It is found that there was a strong agreement in the orders of relative importance for the major pollutants among six major cities in North China, which implies regularities in the meteorological processes of air pollution in that region. In particular, $PM_{2.5}$, PM_{10} , SO_2 and CO shared a common variable importance order and were mostly impacted by the dew point temperature and air pressure, while the NO_2 and O_3 were mostly influenced by the boundary layer height (BLH) and temperature, respectively. We evaluate the impacts of BLH on the pollution levels given the surface meteorological variables. It is found that BLH can be well modeled by the surface meteorological variables. Thus, air quality assessment without using BLH would also produce adequate results.

1. Introduction

Air pollution is an enduring challenge encountered by many countries including China. A substantial part of China has experienced severe air pollution in the last two decades. Epidemiological studies have shown a high correlation between exposure to air pollutants and human health (Donaldson et al., 1998; Schwartz, 2000; Pope III et al., 2002; Chen et al., 2013); research has also found that air pollution may cause social and economic problems (Xie et al., 2016; Zhang et al., 2017b; Feng et al., 2019). Mitigation and control for air pollution are needed in these countries, which requires understanding the mechanisms of the air pollution with respect to its main drivers.

Air pollution is impacted by both the underlying emission and the meteorological condition. As one of the two main drivers, it is important to understand the meteorological aspects of air pollution based on the air quality and meteorological data which are quite available in this era of data. The relationship between meteorological factors and the pollutant concentration has been considered in a range of studies. Tai et al. (2010) studied the correlation of $PM_{2.5}$ and its components with meteorological variables in the United States by applying multiple linear regression and found meteorological variables, can explain up to 50% of the variation of $PM_{2.5}$. Liang et al. (2015, 2016) revealed that $PM_{2.5}$ concentration was highly influenced by meteorological conditions in Beijing and found that 75% of the hourly variation of $PM_{2.5}$ was driven by meteorological factors. Shen et al. (2018) studied the monthly $PM_{2.5}$ levels in Beijing's winters during 2010-2017 and found that 81% of the variation can be explained by the first principal component of 850hPa meridian wind velocity and the relative humidity. Vu et al. (2019) applied the random forest to assess the impact of clean air action on six major pollutants in Beijing by decoupling the impact of meteorology on ambient air quality. Stafoggia et al. (2019) estimated the daily PM_{10} and $PM_{2.5}$ concentrations in Italy during 2013–2015 using the random forest with similar meteorological variables as well as the planetary boundary layer height (BLH). Choubin et al. (2020)

*Corresponding author. Center of Statistical Research and School of Statistics, Southwestern University of Finance and Economics, Chengdu, 611130, China.

**Corresponding author. School of Mathematical Sciences and Guanghua School of Management, Peking University, Beijing, 100871, China.
E-mail address: yaxuan_huang@berkeley.edu (Y. Huang); guobin@swufe.edu.cn (B. Guo); hxsun@pku.edu.cn (H. Sun); liuhuijie16@pku.edu.cn (H. Liu); songxichen@pku.edu.cn (S.X. Chen)
ORCID(s): 0000-0002-2338-0873 (S.X. Chen)

54 studied the hazard prediction of PM_{10} in Barcelona Province, Spain and found seven critical features among a total of
55 thirteen features, which included topographic wetness index and precipitation.

56 In an attempt to construct timely statistical measures to reflect the underlying emission based on the air quality
57 data, Liang et al. (2015) and Zhang et al. (2020) proposed meteorologically adjusted pollution measures (average and
58 quantile) under the meteorological baseline distribution constructed from the readily available surface meteorological
59 data. They showed that the meteorologically adjusted pollution measures can reflect the underlying changes in the
60 emission, while the unadjusted raw average pollution levels are subject to meteorological confounding and cannot
61 objectively reflect the underlying emission.

62 Despite the wide range of studies on the meteorological effects of the air quality motivated from different contexts,
63 there is still a lack of in-depth understanding on the relative importance of the meteorological variables on different air
64 pollutant species. Specifically, given the inter-dependence among the meteorological variables, is there an established
65 order of relative meteorological importance for an air pollutant in a region? An established order of relative importance
66 would imply a common mechanism in the meteorological process of the air pollutant, which would be useful for the
67 modeling and prediction for the air pollutant.

68 This paper tries to answer the above question by analysing air quality data in six major northern China cities. The
69 orders of the variable importance are determined by the ranks of variables being selected in a step-wise forward selec-
70 tion procedure (Hastie et al., 2009) in the regression of pollution concentration on the meteorological variables. The
71 forward variable selection orders take into account the inter-dependence among the meteorological variables relative
72 to the target pollutant. The study finds that there was a strong agreement in the orders of meteorological variables
73 in affecting the pollution concentration among the six major cities in North China. In particular, $PM_{2.5}$, PM_{10} , SO_2
74 and CO shared a common order of meteorological importance led by the dew point temperature, while NO_2 and O_3
75 were mostly influenced by the boundary layer height (BLH) and air temperature, respectively. These suggest there is
76 regularity in the meteorological processes for the six species of air pollutants in North China.

77 The planetary boundary layer height (BLH) defines the vertical dispersion property and is known to be influential to
78 the ground-level concentration of air pollutants. Studies have used BLH as a key factor in the formation and evolution
79 of heavy air pollution (Liu et al., 2013; Quan et al., 2014; Tang et al., 2016; Miao and Liu, 2019). A negative correlation
80 between BLH and pollution concentration from observed data was found in Zhang et al. (2009); Miao et al. (2017); Gui
81 et al. (2019). Xiang et al. (2019) revealed a significant negative correlation between BLH and $PM_{2.5}$ in Beijing during
82 the winter heavy pollution. Miao et al. (2021) found that the relationships between planetary boundary layer and $PM_{2.5}$
83 and O_3 in the summer of Beijing and Shanghai were different with heavy pollution being associated with lower BLH
84 in Beijing, but higher BLH in Shanghai. That higher pollution being associated with higher BLH in Shanghai was due
85 to confounding of transported pollution from warmer inland under westerly wind, which indicates the importance of
86 treating the multivariate meteorological variables collectively.

87 However, unlike the surface meteorological variables which are readily available, there is a delay of around three
88 months in the BLH data assimilated by the European Centre for Medium-Range Weather Forecasts (ECMWF). There
89 are two pending questions regarding the role of BLH in the air pollution processes. One is how much information
90 contained in BLH can be explained by the surface meteorological variables; and the other is the impacts of BLH on the
91 meteorological adjusted average pollution levels. The study finds that, despite a significant negative correlation with
92 the pollutants, BLH was not highly ranked for the five non- NO_2 species after considering the surface variables in the
93 modeling. Furthermore, more than 77% of the variation in the BLH can be explained by the surface variables in the
94 six cities. Even for NO_2 , where BLH was the most important variable, the effects of BLH on NO_2 was less than 1.24%
95 of the meteorologically adjusted average concentration. Therefore, using only the surface meteorological variables for
96 air quality assessment would produce adequate results.

97 The paper is organized as follows. Section 2 outlines the data and the study region. Research methods are outlined
98 in Section 3, which include the nonparametric regression, the forward variable selection and the meteorological ad-
99 justment. Section 4 reports results of the analysis. Additional information, extra tables and figures are provided in the
100 supplementary materials (SM).

101 2. Study Region and Data

102 The study region contains six major cities in the North China, which include two mega-cities of Beijing and Tianjin,
103 and four cities: Shijiazhuang, Jinan, Zhengzhou and Taiyuan, which are the provincial capitals of Hebei, Shandong,
104 Henan and Shanxi provinces, respectively. The region encompassed by Beijing, Tianjin and the four provinces has been

105 the main battle ground for the air pollution mitigation campaign in China; see Figure S1 in SM for the geographical
106 location of the six cities. The time range of the study was from March 1, 2013 to February 28, 2021, which contained 8
107 seasonal years from spring to winter. The unit of our study is season, consisting of spring from March to May, summer
108 from June to August, fall from September to November, and winter from December to February next year.

109 The air quality data consisted of hourly concentrations of air pollutants ($PM_{2.5}$, PM_{10} , SO_2 , NO_2 , CO and O_3) from
110 53 monitoring sites, which are directly managed by the Ministry of Ecology and Environment (MEE). For ozone, we
111 only focused on the period from noon to 7pm (8-hour O_3), when the concentrations tended to be the highest. The 53
112 air quality sites consisted of 11 sites in Beijing and Tianjin, 7 in Shijiazhuang, 8 from Jinan, Zhengzhou and Taiyuan,
113 respectively.

114 The meteorological data contained hourly surface (two meters) meteorological observations of six variables: air
115 temperature (TEMP) and pressure (PRES), dew point temperature (DEWP), wind direction (WD), cumulative wind
116 speed (IWS), and cumulative precipitation (IPREC). The wind directions were grouped into five broad directions:
117 Northwesterly (NW), Northeasterly (NE), Southwesterly (SW), Southeasterly (SE) and CV, where CV was for the
118 calm (wind speed less than 0.5m/s) and variable wind. The cumulative wind speed was the summation of the wind
119 speed since a wind direction was established and was reset to zero when there was a change in the direction. Similarly,
120 the cumulative precipitation referred to the sum of precipitation since the hour when it rained and was reset to zero
121 when there was an hour without precipitation. Data of the above variables were from 12 weather observing stations
122 from China Meteorological Administration (CMA), where 5 of them were located in Beijing, 3 in Tianjin and 1 each in
123 other four cities. We also considered the reanalyzed BLH data from the ERA5 data set of ECMWF, which contained
124 hourly BLH at a grid size of $0.25^\circ \times 0.25^\circ$. We matched each CMA site to the nearest grid of the ERA5 data for BLH as
125 well as to the nearest air quality site. The time range of surface meteorological and BLH data was from March 1, 2011
126 to February 28, 2021. Here, two more years of meteorological data were employed to construct the meteorological
127 baseline in the meteorological adjustment outlined in Section 3.3. The locations of the 53 air quality sites and the
128 meteorological stations are displayed in Figure S2 in SM.

129 In additional to the reanalyzed BLH data from ECMWF, there are other approaches to estimate BLH, which include
130 the radiosonde approach by NOAA (Durre and Yin, 2008) and China's CMA (Guo et al., 2016). The BLH from the
131 radiosonde observations have usually two (sometimes four in summer) observations per day. Seidel et al. (2012)
132 analysed BLH measurements from four approaches over the United States and Europe, and found it was suitable
133 to estimate BLH based on the bulk Richardson number (Vogelezang and Holtslag, 1996), which is the commonly
134 used approach adopted by both the radiosonde observation and the reanalysis approach that produced the ECMWF's
135 ERA5 reanalysis BLH dataset. Based on the radiosonde data of CMA, the properties of planetary boundary layer
136 and its relationship with other meteorological variables in China have been studied. Guo et al. (2019) observed a
137 shift in the temporal trend of BLH in China and found BLH was negatively associated with relative humidity, while
138 positively associated with the near-surface temperature. Zhang et al. (2018) studied the thermodynamic stability of
139 planetary boundary layer in China in summer. They found that convective boundary layer (CBL) dominated in summer
140 throughout China (70%), with sometimes neutral boundary layer (NBL, 26%) and stable boundary layer (SBL, 4%)
141 and BLH of CBL and NBL was positively (negatively) associated with near-surface temperature (humidity), whereas
142 no apparent relationship was found for SBL.

143 Studies (Seidel et al., 2012; Guo et al., 2016) had shown that the BLH estimates from the radiosonde and the
144 ERA Interim, the predecessor of the ERA5, had consistent pattern. Seidel et al. (2012) evaluated the BLH data at
145 the United States and Europe and found that the radiosonde and reanalysis datasets show similar patterns of spatial
146 and seasonal variability though with biases that vary spatially, seasonally and diurnally. Guo et al. (2016) compared
147 the BLH measurements from the radiosonde measurement from China's CMA and the ERA Interim and found good
148 agreement between the two sources. Another aspect is that the BLH data from the ERA5 is at the hourly frequency
149 which matches the frequency of the air quality data. Therefore, we choose the BLH reanalysis data from the ERA5 of
150 ECMWF in our analysis.

151 3. Methods

152 We introduce the regression model used in the study, the forward variable selection that determines the order of
153 the meteorological variables' importance and the meteorological adjustment approach.

154 3.1. Nonparametric Regression Models

155 To avoid potential model mis-specification, we consider nonparametric regression modeling of each air pollutant
 156 with respect to the meteorological variables. The nonparametric aspect of the model makes it more adaptive to the
 157 complex meteorological processes than the linear or nonlinear parametric regression models.

Let $Y_{ijt}(s)$ denote the concentration of a pollutant, say $PM_{2.5}$, at hour t in season j and year i of a monitoring site s , $\mathbf{X}_{ijt}(s)$ be a d -dimensional vector collecting d meteorological variables, and $U_{ijt}(s)$ be the underlying emission leading to the pollution. The following nonparametric model quantifies the relationship between the concentration $Y_{ijt}(s)$, the emission $U_{ijt}(s)$ and the meteorological $\mathbf{X}_{ijt}(s)$:

$$Y_{ijt}(s) = \tilde{m}_j(U_{ijt}(s), \mathbf{X}_{ijt}(s)) + e_{ijt}(s), \quad t = 1, 2, \dots, n_{ij}, \quad (3.1)$$

158 where $\tilde{m}_j(u(s), \mathbf{x}(s)) = E(Y_{ijt}(s)|U_{ijt}(s) = u(s), \mathbf{X}_{ijt}(s) = \mathbf{x}(s))$ is the regression function, $e_{ijt}(s)$ is the residual
 159 satisfying $E(e_{ijt}(s)|U_{ijt}(s), \mathbf{X}_{ijt}(s)) = 0$ and has the finite second moment, and n_{ij} is the total number of observations
 160 in season j of year i at the site s .

As the emission information is largely not available at the hourly frequency, we have to condition on the observable meteorological variable $\mathbf{X}_{ijt}(s)$ to attain the following observable model

$$Y_{ijt}(s) = m_{ij}(\mathbf{X}_{ijt}(s)) + \epsilon_{ijt}(s), \quad t = 1, 2, \dots, n_{ij}, \quad (3.2)$$

161 where $m_{ij}(\mathbf{x}(s)) = E(Y_{ijt}(s)|\mathbf{X}_{ijt}(s) = \mathbf{x}(s))$ is the regression function by taking the conditional expectation on the
 162 observable meteorological data only, $\epsilon_{ijt}(s)$ is the residual satisfying $E(\epsilon_{ijt}(s)|\mathbf{X}_{ijt}(s)) = 0$, $\text{var}(\epsilon_{ijt}(s)|\mathbf{X}_{ijt}(s)) =$
 163 $\sigma_{ij}^2(\mathbf{X}_{ijt}(s))$. Despite $U_{ijt}(s)$ is latent, its information has been reflected in the regression function $m_{ij}(\cdot)$ via the yearly
 164 index i which was not appeared in Model (3.1). See Liang et al. (2015) and Zhang et al. (2020) for more details on the
 165 linkage between Models (3.1) and (3.2).

In Model (3.2), specific parametric form of the regression function $m_{ij}(\cdot)$ is not assumed to allow model flexibility and non-linear pattern in the pollution generation processes. The yearly seasonal regression function $m_{ij}(\cdot)$ can be estimated by the nonparametric kernel estimator (Härdle, 1990). Specifically, let $k(\cdot)$ be a univariate kernel function which is a symmetric probability density function, for instance the Gaussian kernel $k(u) = (2\pi)^{-1/2} \exp(-u^2/2)$, and define the multivariate product kernel

$$K_{\mathbf{h}}(\mathbf{x}) = (h_1 h_2 \dots h_d)^{-1} k(x_1/h_1) k(x_2/h_2) \dots k(x_d/h_d),$$

where $\mathbf{x} = (x_1, x_2, \dots, x_d)'$ and $\mathbf{h} = (h_1, h_2, \dots, h_d)'$ is a vector of smoothing bandwidths. The kernel estimator of $m_{ij}(\mathbf{x})$ is

$$\hat{m}_{ij}(\mathbf{x}) = \frac{\sum_{t=1}^{n_{ij}} K_{\mathbf{h}}(\mathbf{x} - \mathbf{X}_{ijt}(s)) Y_{ijt}(s)}{\sum_{t=1}^{n_{ij}} K_{\mathbf{h}}(\mathbf{x} - \mathbf{X}_{ijt}(s))}. \quad (3.3)$$

166 The smoothing bandwidths \mathbf{h} can be selected by the cross validation method (Härdle, 1990; Liang et al., 2015).

167 3.2. Forward Variable Selection

168 An important question in the air quality assessment is the relative importance of the meteorological variables on the
 169 average pollution level through the regression function $m_{ij}(\cdot)$. We consider the forward variable selection procedure
 170 for the nonparametric regression model, which selects the most influential variable once at a time that achieves the
 171 best fitting performance (Hastie et al., 2009).

To find the most important (the first one to be selected) variable for a pollutant, we start with the univariate model

$$Y_{ijt}(s) = m_{ij}^{(1)}(X_{ijt}^{(1)}(s)) + \epsilon_{ijt}^{(1)}(s), \quad t = 1, 2, \dots, n_{ij}, \quad (3.4)$$

172 where $X_{ijt}^{(1)}(s)$ is a candidate covariate, say DEWP or BLH, $m_{ij}^{(1)}(X_{ijt}^{(1)}(s)) = E(Y_{ijt}(s)|X_{ijt}^{(1)}(s))$ is the conditional mean
 173 given the covariate and $\epsilon_{ijt}^{(1)}(s)$ is the residual. It is noted here and throughout this subsection that the function $m_{ij}^{(1)}(\cdot)$
 174 changes with respect to different covariate $X_{ijt}^{(1)}(s)$ and target cities.

Let $\hat{m}_{ij}^{(1)}(\cdot)$ be the kernel estimate of $m_{ij}^{(1)}(\cdot)$ according to the generic estimator (3.3). The fitting Mean Square Error (MSE) of each variable $X_{ijt}^{(1)}(s)$ at site s for year i and season j is

$$\text{MSE}_{ij}^{(1)}(X_{ijt}^{(1)}(s), s) = \frac{1}{n_{ij}} \sum_{t=1}^{n_{ij}} \{Y_{ijt}(s) - \hat{m}_{ij}^{(1)}(X_{ijt}^{(1)}(s))\}^2. \quad (3.5)$$

To evaluate the explanation power of the meteorological variables, we use the coefficient of determination R^2 to assess fitting performance of the covariate. For nonparametric regression, Doksum and Samarov (1995) proposed to use

$$R_{ij}^2(1, s) = \frac{[\sum_{t=1}^{n_{ij}} \{Y_{ijt}(s) - \bar{Y}_{ij}(s)\} \{\hat{m}_{ij}^{(1)}(X_{ijt}^{(1)}(s)) - \bar{Y}_{ij}(s)\}]^2}{\sum_{t=1}^{n_{ij}} \{Y_{ijt}(s) - \bar{Y}_{ij}(s)\}^2 \sum_{t=1}^{n_{ij}} \{\hat{m}_{ij}^{(1)}(X_{ijt}^{(1)}(s)) - \bar{Y}_{ij}(s)\}^2}, \quad (3.6)$$

175 where $\bar{Y}_{ij}(s) = \sum_{t=1}^{n_{ij}} Y_{ijt}(s)/n_{ij}$ is the average of $Y_{ijt}(s)$ in season j of year i .

176 We use the MSE to evaluate the importance of explanatory variables. By taking the average of $\text{MSE}_{ij}^{(1)}(X_{ijt}^{(1)}(s), s)$
 177 for all the monitoring sites in the city, we obtain the MSE of the city for using the univariate regressor $X_{ijt}^{(1)}$ at year i
 178 and season j . The most important explanatory variable for the city, denoted as $X_{ijt}^{(1)*}$, at year i and season j is the one
 179 that produced the smallest MSE of the city.

To select the second most important regressor, we consider the bivariate regression model

$$Y_{ijt}(s) = m_{ij}^{(2)}(\mathbf{X}_{ijt}^{(2)}(s)) + \epsilon_{ijt}^{(2)}(s), \quad t = 1, 2, \dots, n_{ij}, \quad (3.7)$$

where $\mathbf{X}_{ijt}^{(2)}(s) = (X_{ijt}^{(1)*}(s), X_{ijt}^{(2)}(s))'$ for a $X_{ijt}^{(2)}(s)$ other than the already selected $X_{ijt}^{(1)*}(s)$, $m_{ij}^{(2)}(\mathbf{X}_{ijt}^{(2)}(s)) = E(Y_{ijt}(s)|\mathbf{X}_{ijt}^{(2)}(s))$
 is the conditional mean given $\mathbf{X}_{ijt}^{(2)}(s)$ and $\epsilon_{ijt}^{(2)}(s)$ is the corresponding residual. Similar to selecting the first variable,
 the MSE with $X_{ijt}^{(2)}(s)$ as covariates is

$$\text{MSE}_{ij}^{(2)}(\mathbf{X}_{ijt}^{(2)}(s), s) = \frac{1}{n_{ij}} \sum_{t=1}^{n_{ij}} \{Y_{ijt}(s) - \hat{m}_{ij}^{(2)}(\mathbf{X}_{ijt}^{(2)}(s))\}^2, \quad (3.8)$$

where $\hat{m}_{ij}^{(2)}(\cdot)$ is also attained via (3.3). We select the second important variable by minimizing the average MSE in
 (3.8) over the sites in the city. The coefficient of determination $R^2(2, s)$ by using $\mathbf{X}_{ijt}^{(2)}(s)$ is

$$R_{ij}^2(2, s) = \frac{[\sum_{t=1}^{n_{ij}} \{Y_{ijt}(s) - \bar{Y}_{ij}(s)\} \{\hat{m}_{ij}^{(2)}(\mathbf{X}_{ijt}^{(2)}(s)) - \bar{Y}_{ij}(s)\}]^2}{\sum_{t=1}^{n_{ij}} \{Y_{ijt}(s) - \bar{Y}_{ij}(s)\}^2 \sum_{t=1}^{n_{ij}} \{\hat{m}_{ij}^{(2)}(\mathbf{X}_{ijt}^{(2)}(s)) - \bar{Y}_{ij}(s)\}^2}. \quad (3.9)$$

180 It is shown in Theorem S1 of SM that the variance of $\epsilon_{ijt}^{(2)}(s)$ in Model (3.7) is less than or equal to the variance
 181 of $\epsilon_{ijt}^{(1)}(s)$ in Model (3.4) with the $X_{ijt}^{(1)*}(s)$ as the regressor. This informs the benefits of adding more regressors to the
 182 regression model as it can reduce the variations of the residual term and hence makes the regression model having better
 183 fitting performance, although it is another matter for prediction as having more variables may not lead to better forecast.
 184 An implication of the result is that asymptotically $\text{MSE}_{ij}^{(2)}(\mathbf{X}_{ijt}^{(2)}(s), s) \leq \text{MSE}_{ij}^{(1)}(X_{ijt}^{(1)*}(s), s)$ and the R^2 monotonically
 185 increases with the nested regressors.

186 We proceed the nested forward procedure to select the third, the fourth important covariates and beyond until all
 187 the candidate covariates are considered. The order of the variables being selected contains useful information on the
 188 relative importance of the meteorological factors on the pollutant concentration and has taken into consideration of the
 189 inter-dependence of the meteorological variables.

190 To find out how much information in the BLH that can be explained by the surface meteorological variables, we
 191 carry on the same nested forward selection procedure by substituting the pollution concentration with BLH as the
 192 response variable and the surface variables as the regressors. The order of the variables being selected are used as the
 193 order of relative importance for BLH.

Furthermore, to avoid over-fitting of the nonparametric regression, we used a 10-fold cross validation (CV) procedure to gain information on the out-of-sample performance with the selected variables. The 10-fold CV randomly divides the data into 10 segments and alternatively uses any 9 segments at a time to fit the regression model and the remaining one to gain validation performance of the model. As the validation is made on the data set which is not used in the model fitting, it is more objective and is called the out-of-sample validation. Specifically, we randomly divided the data of a season at a station into 10 segments. Then, for each data segment (as the validation data set), we fit the nonparametric regression model on the remaining 9 segments of data (the training set) and evaluate its performance on the validation data set by calculating the out-of-sample $R_{CV,k}^2$

$$R_{CV,k}^2 = \frac{\{\sum_{t=1}^{n_k} (Y_{k,t} - \bar{Y}_k)(\hat{m}_{-k}(\mathbf{X}_{k,t}) - \bar{Y}_k)\}^2}{\sum_{t=1}^{n_k} (Y_{k,t} - \bar{Y}_k)^2 \sum_{t=1}^{n_k} \{\hat{m}_{-k}(\mathbf{X}_{k,t}) - \bar{Y}_k\}^2}, \quad k = 1, 2, \dots, 10, \quad (3.10)$$

where $\hat{m}_{-k}(\cdot)$ is the fitted regression function without the k th segment, $\mathbf{X}_{k,t}$ and $Y_{k,t}$ denote the t th observation in the k th segment. The cross-validated coefficient of determination R_{CV}^2 is

$$R_{CV}^2 = \frac{1}{10} \sum_{k=1}^{10} R_{CV,k}^2, \quad (3.11)$$

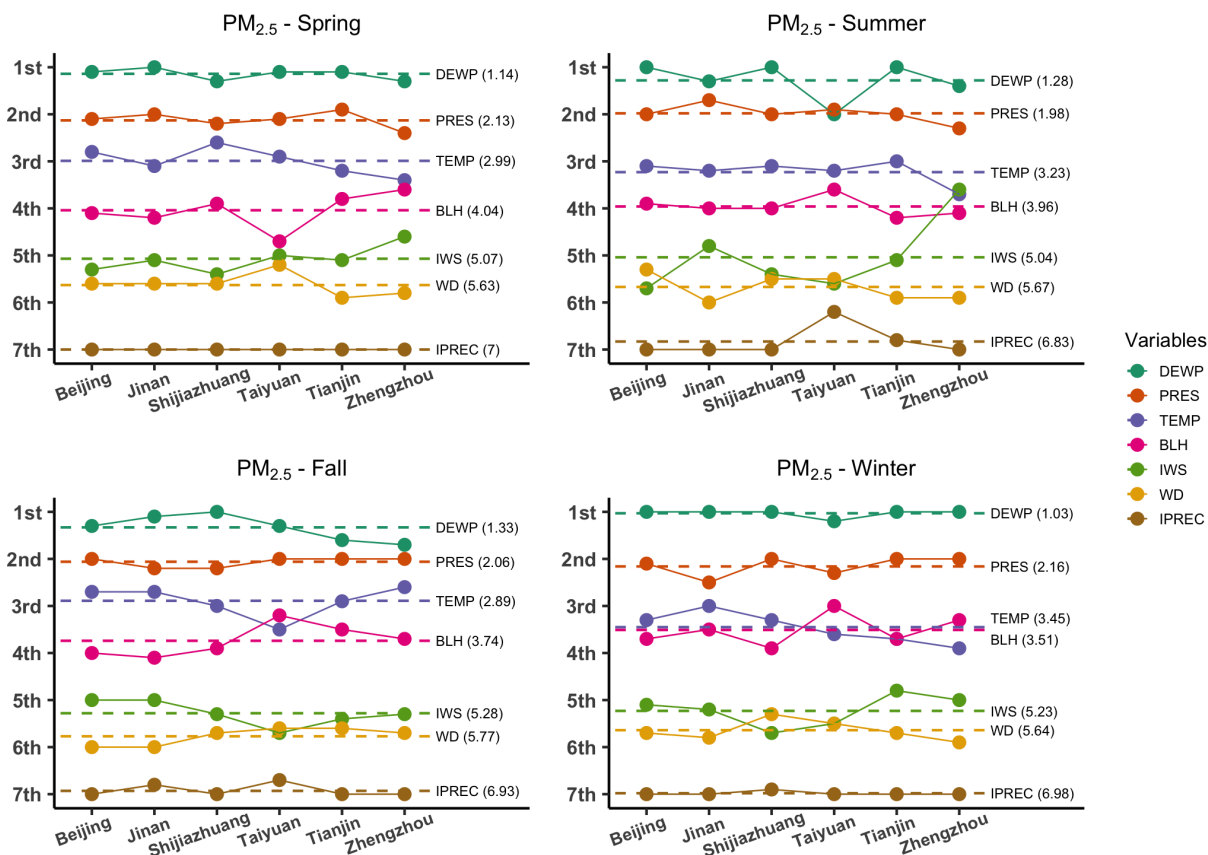


Figure 1: Seasonal average ranks of the meteorological variables for PM_{2.5} regression from 2013 to 2020. The dots with different colors represent ranks with the horizontal dashed lines marking the average ranks, reported inside the parentheses on the right edges, of the variables among the six cities based on the seasonal regression models.

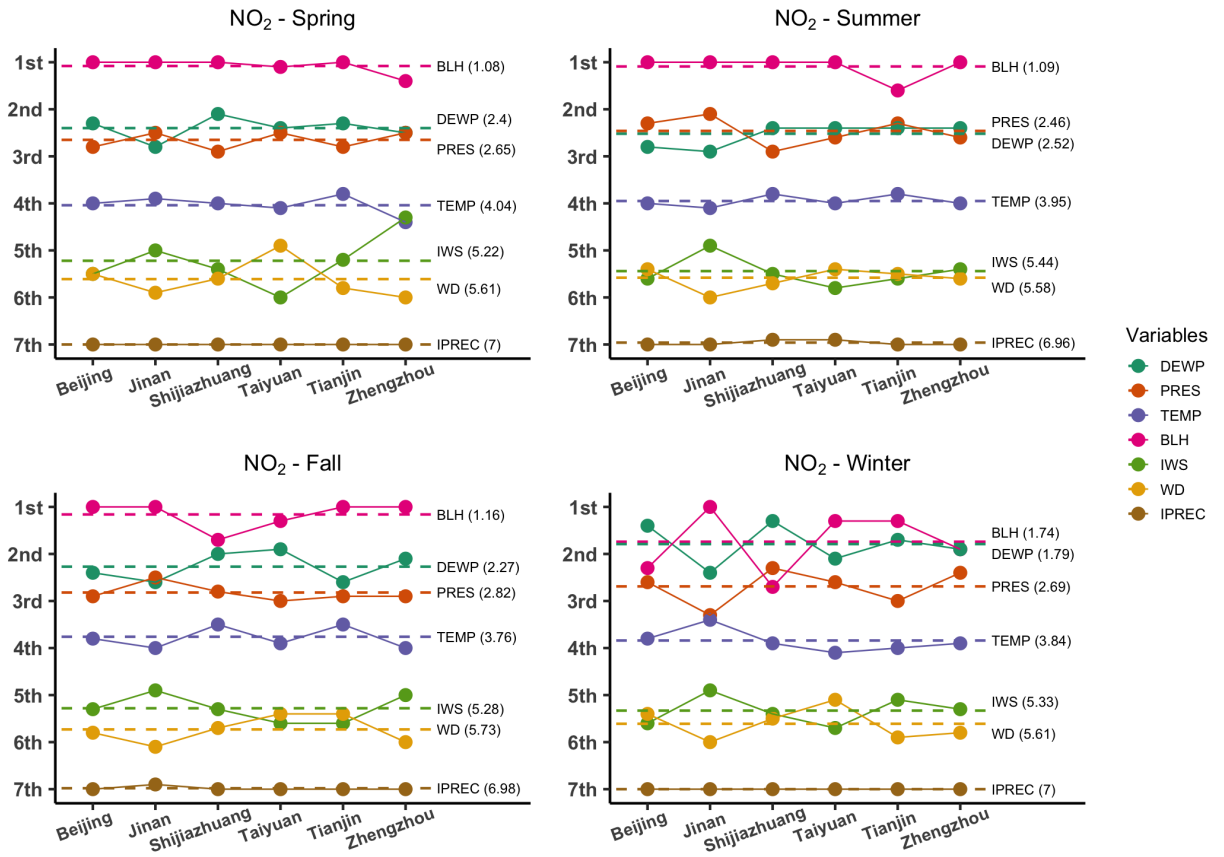


Figure 2: Seasonal average ranks of the meteorological variables for NO₂ regression from 2013 to 2020. The dots with different colors represent ranks with the horizontal dashed lines marking the average ranks, reported inside the parentheses on the right edges, of the variables among the six cities based on the seasonal regression models.

3.3. Meteorological Adjustment

To gain information on the underlying emission, the meteorological effects on the pollution concentration has to be removed. There were methods to adjust for the meteorological variation, including the trend analysis approach proposed in Thompson et al. (2001) based on the linear regression and the three-year moving average method proposed by the US Environmental Protection Agency (EPA). Chen et al. (2018) showed that the latter method can not remove the meteorological confounding. Liang et al. (2015) proposed a general framework to adjust for the meteorological confounding, see Zhang et al. (2020) for comprehensive theoretical analysis.

For air quality evaluation based on the meteorological variables $\mathbf{X}_{ijt}(s)$, we consider Model (3.2) to quantify the meteorological relationship with the pollutant's concentration $Y_{ijt}(s)$. To remove the meteorological confounding, Liang et al. (2015) proposed a version of the average pollution concentration

$$\mu_{ij}(s) = E(Y_{ijt}(s)) = E\{m_{ij}(\mathbf{X}_{ijt}(s))\} = \int m_{ij}(\mathbf{x}, s) f_j(\mathbf{x}, s) d\mathbf{x}, \quad (3.12)$$

where $m_{ij}(\mathbf{X}_{ijt}(s))$ is the regression function that can be estimated by (3.3), and $f_j(\mathbf{x}, s)$ is a baseline probability density function for the meteorological variable $\mathbf{X}(s)$ in season j which can be constructed as

$$f_j(\mathbf{x}, s) = A_j^{-1} \sum_{a=1}^{A_j} f_{aj}(\mathbf{x}, s),$$

where $f_{aj}(\mathbf{x}, s)$ is the density function of meteorological variable $\mathbf{X}_{ajt}(s)$ in year a of season j , and A_j is the number

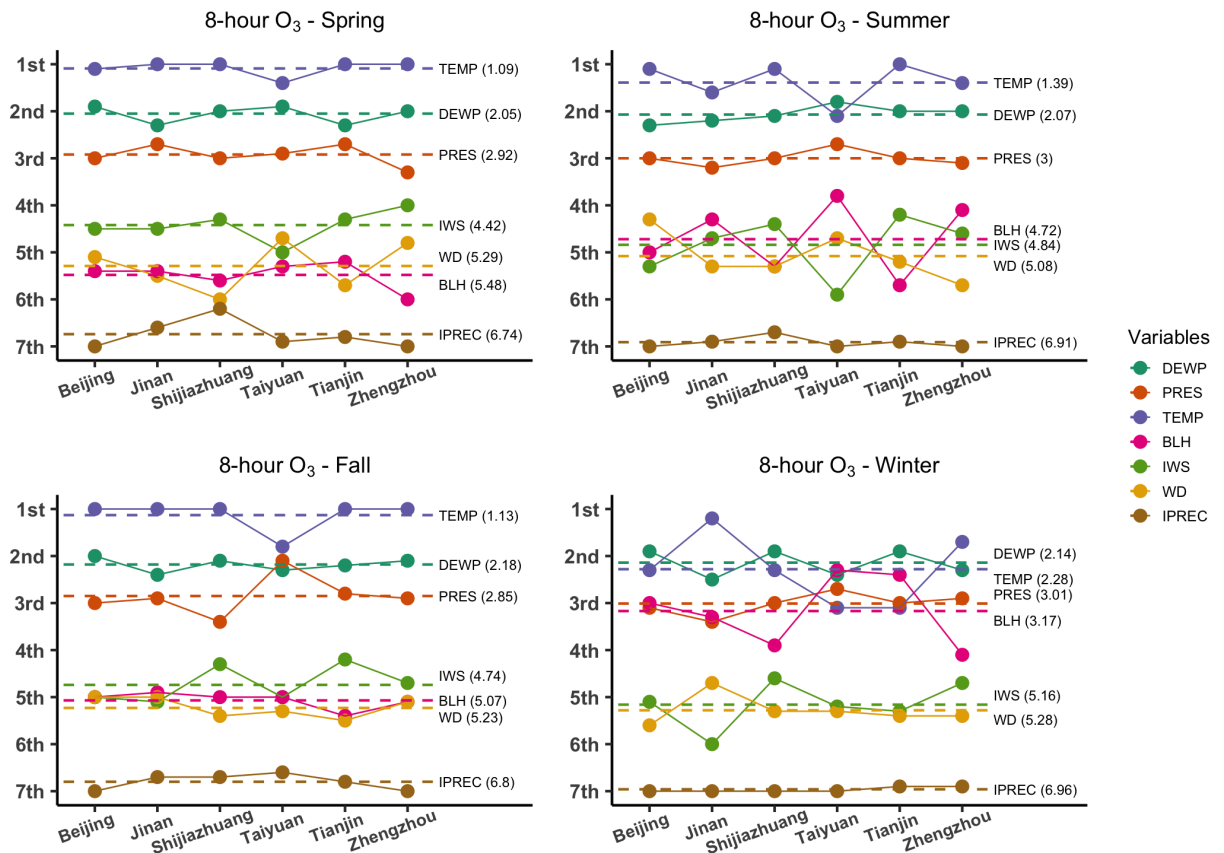


Figure 3: Seasonal average ranks of the meteorological variables for 8-hour O₃ regression from 2013 to 2020. The dots with different colors represent ranks with the horizontal dashed lines marking the average ranks, reported inside the parentheses on the right edges of the variables among the six cities based on the seasonal regression models.

of years that the meteorological data are available. In our study, A_j was 10 as we used 10 years' meteorological data from March 2011 to February 2021 to build the baseline density $f_j(\mathbf{x}, s)$.

An estimator for $\mu_{ij}(s)$ which is free of the meteorological variation and can reflect the emission is

$$\hat{\mu}_{ij}(s) = \left(\sum_{a=1}^{A_j} n_{aj} \right)^{-1} \sum_{a=1}^{A_j} \sum_{t=1}^{n_{aj}} \hat{m}_{ij}(\mathbf{X}_{ajt}(s)), \quad (3.13)$$

where n_{aj} is the number of meteorological sample size in season j of year a . The nonparametric regression estimator $\hat{m}_{ij}(\cdot)$ is (3.3). Unlike the raw averages, the temporally adjusted averages $\hat{\mu}_{ij}(s)$ are comparable for the specific season over different years as shown in Zhang et al. (2017a) and Chen et al. (2018).

The air quality measures $\hat{\mu}_{ij}(\mathcal{A})$ over a region \mathcal{A} occupied by a city is an average of the adjusted $\hat{\mu}_{ij}(s)$ over the monitoring sites in the region:

$$\hat{\mu}_{ij}(\mathcal{A}) = |\mathcal{A}|^{-1} \sum_{s \in \mathcal{A}} \hat{\mu}_{ij}(s), \quad (3.14)$$

where $|\mathcal{A}|$ denotes the number of monitoring sites in the region \mathcal{A} .

4. Results

We present the results of the analyzes in this section.

4.1. Variable Importance for Pollutants

We used the proposed methods in Section 3.2 to obtain the orders of meteorological variables for the six pollutants in each city in each season. Specifically, for the pollutant at each season and year, we attained the order of variable importance for each station in a city and then average the ranks over all stations in the city and over all years for the season.

Figures 1-3 display the seasonal average orders of the meteorological variables for PM_{2.5}, NO₂ and 8-hour O₃ in each city respectively, while those for the other three pollutants are shown in Figures S3-S5. The numerical average ranks for each pollutant in the six cities were reported in Tables S1-S6.

For PM_{2.5}, as shown in Figure 1, there was a consistent ordering in the variable ranks among the seasons and the cities, namely

$$\text{DEWP, PRES, TEMP, BLH, IWS, WD, IPREC.} \quad (4.1)$$

We also noticed that, the above order of variable importance for PM_{2.5} was also obeyed for PM₁₀, SO₂ and CO as shown in Figures S3-S5 in the SM. However, the orders for NO₂ and 8-hour O₃, as shown in Figures 2 and 3 were different. For NO₂, the variable order was

$$\text{BLH, DEWP, PRES, TEMP, IWS, WD, IPREC,} \quad (4.2)$$

which saw BLH jumped to the leading position while the relative order of the other variables remained unchanged from (4.1). For ozone, the order was

$$\text{TEMP, DEWP, PRES, BLH, IWS, WD, IPREC,} \quad (4.3)$$

which had the temperature (TEMP) became the leading variable. The strong consistency in the three orderings (4.1)-(4.3) among the six cities indicated consistency of the meteorological effects on the three categories of pollutants in the study region.

The dew point temperature (DEWP) was the highest ranked variable for the PM-CO-SO₂ group, and the second ranked for NO₂ and O₃. The reason for DEWP's leading roles was its ability in reflecting both temperature and humidity. It is noted that we did not consider the relative humidity as a variable since it is determined by the dew point and the air temperature. Air pressure (PRES) was the second ranked for the PM-CO-SO₂ group, and the third ranked for NO₂ and O₃. Air temperature (TEMP) was the third ranked for the PM-CO-SO₂ group, the fourth for NO₂ but risen to the top place for ozone. TEMP's leading role for ozone was due to its being a proxy for the solar radiation, a major element in the photo-chemistry process of ozone generation. The wind direction (WD) and the cumulative precipitation (IPREC) ranked the lowest for all six pollutants. Cumulative wind speed (IWS) was the third last ranked for all six species.

It should be noted that the ranking did not suggest that the lower ranking variables, say IWS, WD and IPREC, were not influential to the pollutant's concentration, but rather their relative importance after considering the other variables. In particular, the impacts of the lowly ranked variables had been represented by the higher ranking variables. Indeed, a change in the wind direction from the southerly to northerly would bring dryer, cooler and cleaner air mass from the north, which led to a lower DEWP.

It was interesting to see BLH became the leading variable for NO₂. BLH was the clear top variable for all six cities in the three non-winter seasons. The leading role of BLH to NO₂ can be understood by noting the pathway $\text{NO}_2 + \text{O}_2 \rightarrow \text{O}_3 + \text{NO}$ that converts NO₂ to ozone under solar radiation, and the fact that the reverse pathway is weakened by interference of other chemical elements in polluted air which consumed much NO leading to a weakening of the reverse process. As BLH is highly correlated with the solar radiation, BLH became a proxy for the radiation and hence the leading meteorological variable for NO₂. Indeed, as shown in Figures S6 and S7, NO₂ was the pollutant which BLH had both the highest correlation and the partial correlation after conditioning on the other meteorological variables. The partial correlation has filtered out the effects of the other variables on the pollutant (Anderson, 2003). While the partial correlation of BLH with the other five pollutants were much reduced after considering the other variables, its partial correlation with the NO₂ remained at a relative high level as shown in Figures S12 and S13, which provides a statistical understanding for the BLH's leading role in NO₂.

The temperature's leading role to Ozone can be appreciated by its close relationship with the solar radiation, and the latter is directly involved in the ozone generation. As shown in Figures S10 and S11, the temperature had the highest

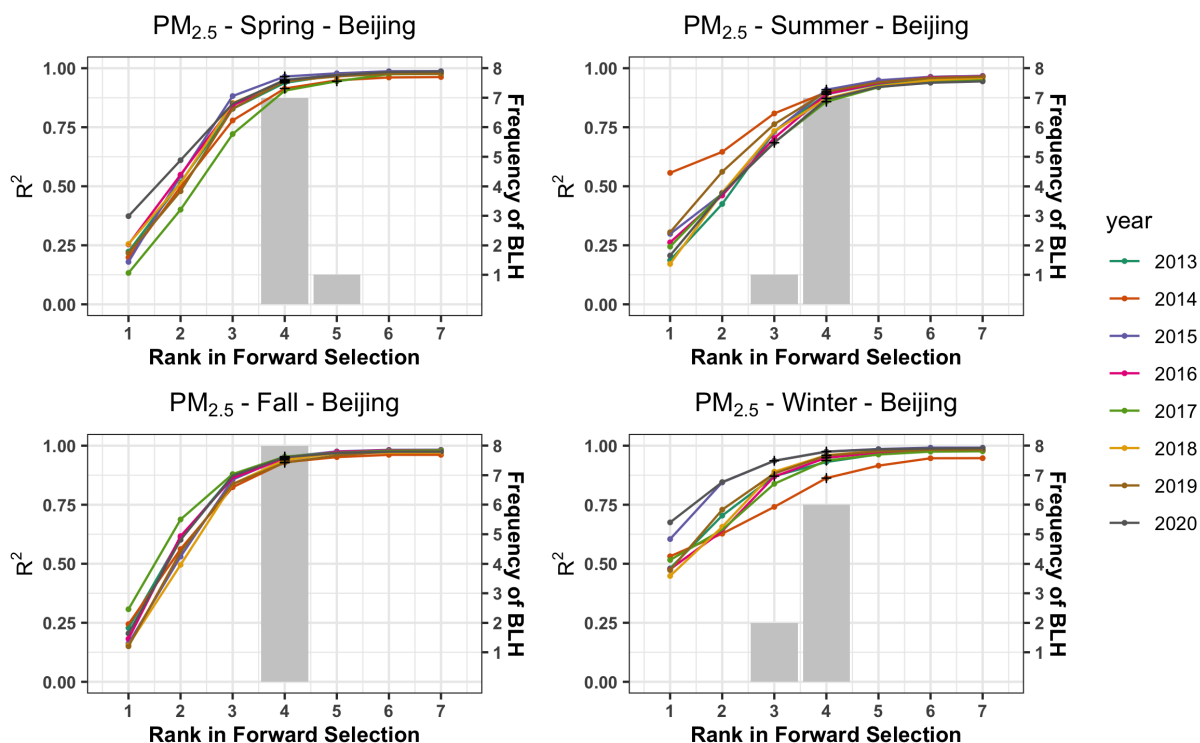


Figure 4: Seasonal step-wise cumulative in-sample R^2 (segmented lines) with variables added to the nonparametric regression of $PM_{2.5}$ in Beijing according to their ranks selected in the forward procedure, and the histograms (grey color bars) for the distributions of the ranks occupied by the BLH (scaled on the right). The ranks of BLH are also marked by black crosses.

249 partial correlation with the O_3 which explained its leading role among the meteorological variables. It is interesting
 250 to see a large change in the correlation between BLH and ozone after we conditioned on the other variables in Figures
 251 S6 and S7, which provided an explanation for BLH's dropping to the fifth or sixth place for ozone in the non-winter
 252 seasons in Figure 3.

253 Figure 4 reports the cumulative R^2 curves in Beijing as the meteorological variables were added according to the
 254 forward selection procedure. Similar plots for the other five cities are available in Figures S14-S16. These figures
 255 showed a rapid increase in the R^2 as the first three or four ranked variables were added in the models, while the growth
 256 in the cumulative R^2 was much slower after the fourth variables. These suggest that the top three or four variables in
 257 (4.1)-(4.3) contain high percentages of the meteorological information.

258 4.2. Role of boundary layer height

259 The analysis in the last subsection showed that, except for NO_2 , the BLH ranked around the fourth or fifth place in
 260 the overall variable ranking for the pollutants. In particular, it ranked behind DEWP, PRES, TEMP for the five non-
 261 NO_2 pollutants in both the ranking and the contributions to the R^2 in Figure 4. These did not necessarily imply that
 262 BLH was not influential to these species, but rather its role was relatively less important when the other meteorological
 263 variables were considered.

264 For better viewing on the correlation between BLH and the pollutants, Figures S6 and S7 display the heat-map
 265 of the correlations between BLH and the pollutants and the corresponding partial correlations by conditioning on the
 266 surface meteorological variables. The partial correlations filter out the role of the surface meteorological variables,
 267 and tended to reduce the magnitude of the usual correlations. Both forms of the correlations were largely negative for
 268 the five pollutants other than the 8-hour O_3 which were positive. The negative correlations between BLH and NO_2
 269 were the highest among the six species. This was consistent with the finding in (4.2) that the BLH was the leading
 270 variable for NO_2 .

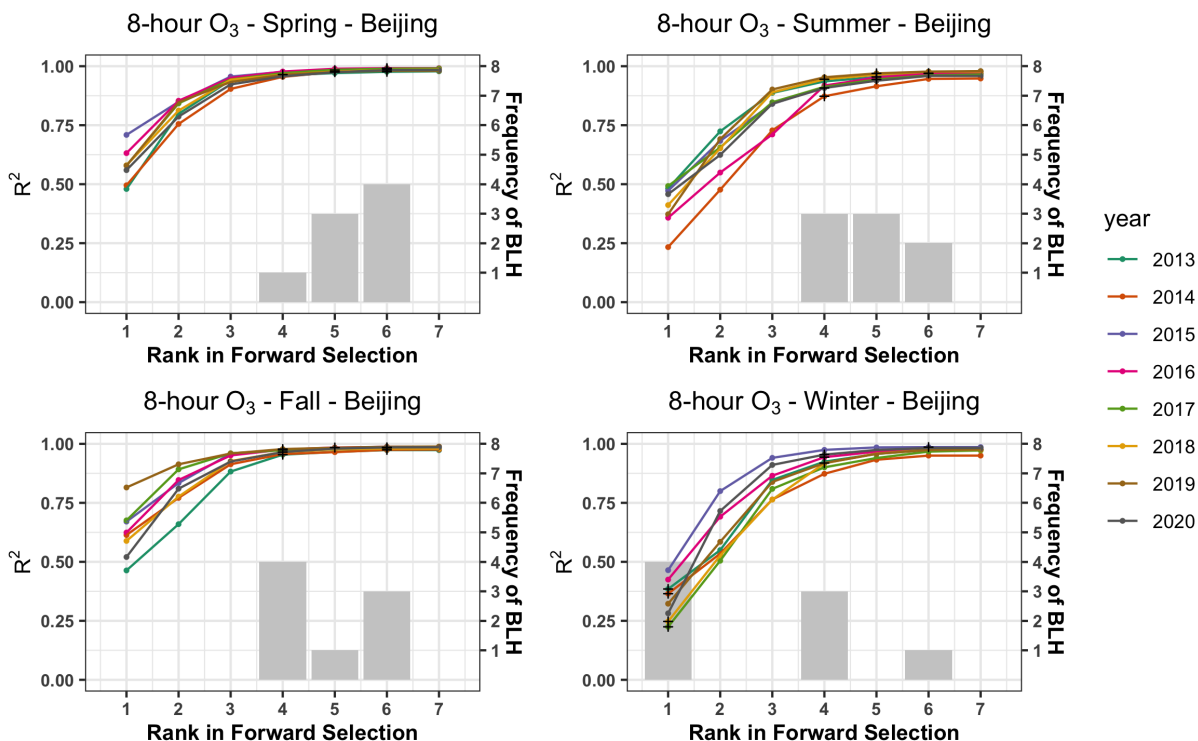


Figure 5: Seasonal step-wise cumulative in-sample R^2 (segmented lines) with variables added to the nonparametric regression of 8-hour O_3 in Beijing according to their ranks selected in the forward procedure, and the histograms (grey color bars) for the distributions of the ranks occupied by the BLH (scaled on the right). The ranks of BLH are also marked by black crosses.

Now we try to understand the importance of BLH from the aspect that how much data information contained in BLH can be explained by the other meteorological variables. To address this question, we conduct analysis based on the nonparametric regression model (3.2) but with the response variable Y being BLH and the regressors $\tilde{\mathbf{X}}$ being the other meteorological variables:

$$BLH_{ijt} = m_{ij}(\tilde{\mathbf{X}}_{ijt}) + \epsilon_{ijt}, \quad t = 1, 2, \dots, n_{ij}. \quad (4.4)$$

271 Figure S23 in SM reports the variable ranks from the forward selection algorithm, which shows that the temperature
 272 (TEMP), dew point (DEWP), and pressure (PRES) were the top three factors for BLH in three seasons other than fall,
 273 where in fall IWS replaced pressure as the third one. We consider two different subsets of $\tilde{\mathbf{X}}$. One was the five variables
 274 without IPREC as it was the last ranked in the rankings reported in (4.1)-(4.3). And the other subset contains only
 275 the three top variables, DEWP, TEMP and PRES, which were the most important three ones among most cities and
 276 seasons and they ranked in front of BLH in most cases.

277 Table 1 reports both the seasonal in-sample R^2 and the out-of-sample R^2_{CV} when regressing BLH with respect to
 278 the three and the five variable sets, respectively, for the six cities. The results showed that the average in-sample R^2
 279 with the three top variables can attain at least 77% of the total variation in BLH among the four seasons and the six
 280 cities, which was risen to at least 92% with the five variables' set of the meteorological factors as the regressors. As
 281 expected, the out-of-sample R^2_{CV} s were lower than the corresponding R^2 . However, the minimum R^2_{CV} was still 57%
 282 with the three top variables and risen to 68% for the five variables' set.

283 The results conveyed in Table 1 suggest that much of the information in BLH can be modeled by the other meteorological
 284 variables, and by the top three variables DEWP, TEMP and PRES in particular. The results provide explanation
 285 on why BLH was not ranked high for most of the pollutant species as conveyed in Figures 1 and 3, as well as on the
 286 marginal contribution to the R^2 by the BLH (Figures 4 and 5).

Table 1

Seasonal in-sample R^2 and the out-sample R_{CV}^2 for BLH with the top three meteorological variables (DEWP, PRES, TEMP) or five variables (DEWP, PRES, TEMP, IWS, WD) for the six cities (averaged from 2013 to 2020). The numbers inside the parentheses are the standard errors of the average R^2 above.

City	Number of Variables	Spring		Summer		Fall	Winter		
		R^2	R_{CV}^2	R^2	R_{CV}^2	R^2	R_{CV}^2	R^2	R_{CV}^2
Beijing	3	0.85	0.64	0.77	0.57	0.88	0.71	0.78	0.60
	5	0.96	0.76	0.91	0.67	0.95	0.77	0.93	0.69
Jinan	3	0.82	0.56	0.78	0.57	0.86	0.64	0.80	0.56
	5	0.95	0.68	0.94	0.69	0.96	0.76	0.95	0.71
Shijiazhuang	3	0.83	0.59	0.76	0.54	0.85	0.63	0.78	0.54
	5	0.96	0.73	0.92	0.66	0.95	0.73	0.93	0.68
Taiyuan	3	0.78	0.58	0.78	0.63	0.82	0.66	0.76	0.58
	5	0.94	0.71	0.90	0.70	0.93	0.72	0.91	0.66
Tianjin	3	0.83	0.58	0.76	0.55	0.85	0.67	0.79	0.57
	5	0.96	0.73	0.92	0.67	0.96	0.76	0.94	0.71
Zhengzhou	3	0.82	0.58	0.77	0.57	0.85	0.63	0.81	0.58
	5	0.96	0.74	0.94	0.71	0.96	0.74	0.95	0.70
Average	3	0.82 (0.009)	0.59 (0.011)	0.77 (0.003)	0.57 (0.012)	0.85 (0.008)	0.66 (0.012)	0.79 (0.007)	0.57 (0.008)
	5	0.96 (0.003)	0.72 (0.010)	0.92 (0.006)	0.68 (0.009)	0.95 (0.005)	0.75 (0.008)	0.93 (0.006)	0.69 (0.008)

To gain further insights on the role of the BLH, we studied its influence on the air quality measures by calculating the meteorologically adjusted average concentration (3.13) with BLH versus those without BLH as a covariate. Specifically, for a city that encompasses region \mathcal{A} , let $\hat{\mu}_{ij}(\mathcal{A})$ and $\hat{\mu}_{ij}^{-BLH}(\mathcal{A})$ be the adjusted average concentrations with and without BLH, respectively, for year i and season j via (3.14). The BLH effect is measured by relative change caused by without BLH

$$\hat{r}_{ij}(\mathcal{A}) = |\hat{\mu}_{ij}(\mathcal{A}) - \hat{\mu}_{ij}^{-BLH}(\mathcal{A})| / \hat{\mu}_{ij}(\mathcal{A}). \quad (4.5)$$

287 Table 2 summarizes the average and the maximum BLH effects (in percentage) over the 32 seasons from 2013 to
 288 2020 for the six pollutants in the six cities. Among the 36 pollutant-city combinations for the estimated average BLH
 289 effects, 16 of them were less than 1% and 15 between 1% and 1.5%, and only 5 of them above 1.5%, indicating rather
 290 mild BLH effects on the pollutants among the six cities. However, the BLH effects appeared not evenly distributed
 291 among the species and cities. Beijing and Tianjin had smaller average and maximum effects while those for Taiyuan
 292 and Jinan were larger. Taiyuan had the largest average and maximum effects among the six cities, which may be due
 293 to its higher elevation (average 800 meters) as the other five cities are all situated in the North China Plain.

294 To gain information on the significance of the BLH effects reported in Table 2, we conducted statistical testing for
 295 the hypothesis $H_0 : \mu_{ij}(\mathcal{A}) = \mu_{ij}^{-BLH}(\mathcal{A})$ versus $H_1 : \mu_{ij}(\mathcal{A}) \neq \mu_{ij}^{-BLH}(\mathcal{A})$ at 5% significance level in each season for
 296 each city and pollutant. The p-value of each testing was obtained based on 300 bootstrap resampling of the test statistics
 297 (Liang et al., 2015; Zhang et al., 2017a). Table 3 summarizes the overall testing results by providing the frequencies
 298 of rejecting the above H_0 among the 32 seasons from 2013 to 2020 for each city and each pollutant. It shows that only
 299 Taiyuan, Jinan and Zhengzhou had significant BLH effects. Taiyuan had the most number of significant differences
 300 for three species (SO_2 , NO_2 and O_3), but still the numbers of significant seasons ranged from 2 to 4 seasons. The

Table 2

The average relative BLH effects (standard deviations) and the maximum BLH effects (the second row) $\hat{r}_{ij}(\mathcal{A})$ in (4.5) in percentage term for the six pollutants and the six cities over the 32 seasons from 2013 to 2020.

Pollutants	Beijing	Jinan	Shijiazhuang	Taiyuan	Tianjin	Zhengzhou	Average
PM _{2.5}	0.83(0.12)	1.29(0.15)	0.82(0.12)	1.57(0.24)	1.02(0.13)	0.73(0.09)	1.04(0.12)
	2.54	4.02	2.74	5.52	2.89	1.98	3.28
PM ₁₀	0.65(0.12)	1.07(0.14)	0.67(0.10)	1.34(0.22)	0.93(0.11)	0.70(0.08)	0.89(0.10)
	2.82	4.06	2.27	4.31	2.15	1.80	2.90
SO ₂	1.20(0.13)	1.37(0.19)	1.19(0.21)	2.32(0.36)	1.15(0.13)	0.96(0.14)	1.37(0.18)
	2.74	4.38	6.06	7.90	3.25	3.22	4.59
NO ₂	0.94(0.13)	1.58(0.19)	1.29(0.15)	1.39(0.18)	1.13(0.13)	1.13(0.28)	1.24(0.08)
	2.96	4.29	4.06	3.77	2.63	8.55	4.37
CO	0.78(0.14)	1.24(0.18)	1.10(0.18)	0.99(0.17)	0.80(0.10)	0.72(0.15)	0.94(0.08)
	2.86	4.04	3.71	4.16	2.39	4.37	3.59
8-hour O ₃	0.56(0.08)	2.72(0.38)	1.29(0.21)	3.79(0.47)	0.79(0.11)	0.73(0.12)	1.65(0.49)
	1.76	7.19	5.13	10.62	2.27	2.34	4.88
Average	0.83(0.09)	1.55(0.24)	1.06(0.11)	1.90(0.42)	0.97(0.06)	0.83(0.07)	1.19(0.10)
	2.61	4.66	3.99	6.05	2.60	3.71	3.94

Table 3

Frequencies (percentages) of significant differences (at 5% level) in the seasonal adjusted average pollution concentrations that included BLH versus those without BLH for the six pollutants and the six cities over the 32 seasons from 2013 to 2020.

Pollutants	Beijing	Jinan	Shijiazhuang	Taiyuan	Tianjin	Zhengzhou
PM _{2.5}	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)
PM ₁₀	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)	0 (0)
SO ₂	0 (0)	0 (0)	0 (0)	2 (6.25%)	0 (0)	0 (0)
NO ₂	0 (0)	0 (0)	0 (0)	3 (9.38%)	0 (0)	2 (6.25%)
CO	0 (0)	2 (6.25%)	0 (0)	0 (0)	0 (0)	1 (3.13%)
8-hour O ₃	0 (0)	2 (6.25%)	0 (0)	4 (12.5%)	0 (0)	0 (0)

301 results of Table 3 were consistent to those in Table 2, showing the BLH's effect was largely small as far as the average
 302 air quality measures were concerned.

303 5. Conclusion

304 Existing studies on the meteorological effects on air pollution tended to consider one variable at a time as showed
 305 in the cited works of the paper. This study is designed to evaluate the collective meteorological effects that takes into
 306 account the mutual dependence among the meteorological variables. And it finds that there were much agreement
 307 in the meteorological effects on the air pollution in North China, as reflected by the much agreed orders of variable
 308 importance on the six major air pollutants. These suggest stable meteorological processes with respect to the air

309 pollution in North China, which may be used for meteorological modeling and prediction of air quality.

310 The study also reveals that, due to the inter-dependence of the meteorological variables, the top three variables can
311 provide quite satisfactory modeling for the meteorological effects of a pollutant, as reflected by the rapid increase in the
312 R^2 among the first three variables. The inter-dependence of the meteorological variables means that the less ranked
313 variable's information and their effects on the pollutants can be well represented by the top ranked meteorological
314 variables, as demonstrated in the case of BLH. The inter-dependence also means that adequate surface air quality
315 measures can be well approximated by utilizing the surface meteorological variables as much of the data information
316 in the mixing layer may be well represented by the surface variables.

317 Our study has considered only one non-surface variable. Other variables at the the vertical pressure layers within
318 the mixing layer can be considered and their importance can be assessed using the method demonstrated in the study.
319 Given the high correlation between the variables on the surface and the vertical layer, it would not be surprising to see
320 the meteorological information in the vertical layer can be well substituted by the surface variables for most of the
321 pollutants.

322 **CRedit authorship contribution statement**

323 **Yaxuan Huang:** Software, Formal analysis, Data Curation, Writing - original draft. **Bin Guo:** Conceptualization,
324 Methodology, Software, Formal analysis, Data Curation, Writing - original draft. **Haoxuan Sun:** Software, Formal
325 analysis, Data Curation, Writing - original draft. **Huijie Liu:** Software, Formal analysis, Data Curation. **Song Xi**
326 **Chen:** Conceptualization, Methodology, Writing - original draft, Supervision, Funding acquisition.

327 **Declaration of competing interest**

328 The authors declare that they have no competing financial interests or personal relationships that could have ap-
329 peared to influence the work reported in this paper.

330 **Acknowledgment**

331 The research was partially supported by National Natural Science Foundation of China Grants 92046021, 12071013,
332 71973005 and 11971390, and LMEQF at Peking University. Bin Guo acknowledges support from Center of Statistical
333 Research at SWUFE and Fundamental Research Funds for the Central Universities JBK1806002.

334 **Appendix A. Supplementary data**

335 The authors declare that supplementary data related to this article are available in the following online reposi-
336 tory. The ERA5 hourly data used in this study were collected from the ECMWF website (<https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels?tab=overview>). The pollution data
337 are in the repository (<https://archive.ics.uci.edu/ml/datasets/Beijing+Multi-Site+Air-Quality+Data>).

339 **References**

- 340 Anderson, T.W., 2003. An Introduction to Multivariate Statistical Analysis, 3rd ed., John Wiley and Sons. New York .
- 341 Chen, L., Guo, B., Huang, J., He, J., Wang, H., Zhang, S., Chen, S.X., 2018. Assessing air-quality in Beijing-Tianjin-Hebei region: The method
342 and mixed tales of PM_{2.5} and O₃. *Atmospheric Environment* 193, 290–301.
- 343 Chen, R., Zhao, Z., Kan, H., 2013. Heavy smog and hospital visits in Beijing, China. *American Journal of Respiratory and Critical Care Medicine*
344 188, 1170–1171.
- 345 Choubin, B., Abdolshahnejad, M., Moradi, E., Querol, X., Mosavi, A., Shamshirband, S., Ghamisi, P., 2020. Spatial hazard assessment of the PM10
346 using machine learning models in Barcelona, Spain. *Science of the Total Environment* 701, 134474.
- 347 Doksum, K., Samarov, A., 1995. Nonparametric estimation of global functionals and a measure of the explanatory power of covariates in regression.
348 *Annals of Statistics* 23, 1443–1473.
- 349 Donaldson, K., Li, X., MacNee, W., 1998. Ultrafine (nanometre) particle mediated lung injury. *Journal of Aerosol Science* 29, 553–560.
- 350 Durre, I., Yin, X., 2008. Enhanced radiosonde data for studies of vertical structure. *Bulletin of the American Meteorological Society* 89, 1257–1262.
- 351 Feng, Z., De Marco, A., Anav, A., Gualtieri, M., Sicard, P., Tian, H., Fornasier, F., Tao, F., Guo, A., Paoletti, E., 2019. Economic losses due to
352 ozone impacts on human health, forest productivity and crop yield across China. *Environment International* 131, 104966.
- 353 Gui, K., Che, H., Wang, Y., Wang, H., Zhang, L., Zhao, H., Zheng, Y., Sun, T., Zhang, X., 2019. Satellite-derived PM_{2.5} concentration trends over
354 Eastern China from 1998 to 2016: Relationships to emissions and meteorological parameters. *Environmental Pollution* 247, 1125–1133.

355 Guo, J., Li, Y., Cohen, J.B., Li, J., Chen, D., Xu, H., Liu, L., Yin, J., Hu, K., Zhai, P., 2019. Shift in the temporal trend of boundary layer height in
356 China using long-term (1979–2016) radiosonde data. *Geophysical Research Letters* 46, 6080–6089.

357 Guo, J., Miao, Y., Zhang, Y., Liu, H., Li, Z., Zhang, W., He, J., Lou, M., Yan, Y., Bian, L., Zhai, P., 2016. The climatology of planetary boundary
358 layer height in China derived from radiosonde and reanalysis data. *Atmospheric Chemistry and Physics* 16, 13309–13319.

359 Härdle, W., 1990. *Applied Nonparametric Regression*. Cambridge University.

360 Hastie, T., Tibshirani, R., Friedman, J., 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed., Springer.

361 Liang, X., Li, S., Zhang, S., Huang, H., Chen, S.X., 2016. PM_{2.5} data reliability, consistency, and air quality assessment in five Chinese cities.
362 *Journal of Geophysical Research: Atmospheres* 121, 10220–10236.

363 Liang, X., Zou, T., Guo, B., Li, S., Zhang, H., Zhang, S., Huang, H., Chen, S.X., 2015. Assessing Beijing's PM_{2.5} pollution: severity, weather
364 impact, APEC and winter heating. *Proceedings of the Royal Society A* 471, 20150257.

365 Liu, X.G., Li, J., Qu, Y., Han, T., Hou, L., Gu, J., Chen, C., Yang, Y., Liu, X., Yang, T., Zhang, Y., Tian, H., Hu, M., 2013. Formation and evolution
366 mechanism of regional haze: a case study in the megacity Beijing, China. *Atmospheric Chemistry and Physics* 13, 4501–4514.

367 Miao, Y., Che, H., Zhang, X., Liu, S., 2021. Relationship between summertime co-occurring PM_{2.5} and O₃ pollution and boundary layer height
368 differs between Beijing and Shanghai, China. *Environmental Pollution* 268, 115775.

369 Miao, Y., Guo, J., Liu, S., Liu, H., Li, Z., Zhang, W., Zhai, P., 2017. Classification of summertime synoptic patterns in Beijing and their associations
370 with boundary layer structure affecting aerosol pollution. *Atmospheric Chemistry and Physics* 17, 3097–3110.

371 Miao, Y., Liu, S., 2019. Linkages between aerosol pollution and planetary boundary layer structure in China. *Science of the Total Environment*
372 650, 288–296.

373 Pope III, C.A., Burnett, R.T., Thun, M.J., Calle, E.E., Krewski, D., Ito, K., Thurston, G.D., 2002. Lung cancer, cardiopulmonary mortality, and
374 long-term exposure to fine particulate air pollution. *Journal of the American Medical Association* 287, 1132–1141.

375 Quan, J., Tie, X., Zhang, Q., Liu, Q., Li, X., Gao, Y., Zhao, D., 2014. Characteristics of heavy aerosol pollution during the 2012–2013 winter in
376 Beijing, China. *Atmospheric Environment* 88, 83–89.

377 Schwartz, J., 2000. The distributed lag between air pollution and daily deaths. *Epidemiology* 11, 320–326.

378 Seidel, D.J., Zhang, Y., Beljaars, A., Golaz, J.C., Jacobson, A.R., Medeiros, B., 2012. Climatology of the planetary boundary layer over the
379 continental United States and Europe. *Journal of Geophysical Research: Atmospheres* 117, 17106.

380 Shen, L., Jacob, D.J., Mickley, L.J., Wang, Y., Zhang, Q., 2018. Insignificant effect of climate change on winter haze pollution in Beijing. *Atmo-
381 spheric Chemistry and Physics* 18, 17489–17496.

382 Stafoggia, M., Bellander, T., Bucci, S., Davoli, M., de Hoogh, K., de' Donato, F., Gariazzo, C., Lyapustin, A., Michelozzi, P., Renzi, M., Scortichini,
383 M., Shtein, A., Viegi, G., Kloog, I., Schwartz, J., 2019. Estimation of daily PM₁₀ and PM_{2.5} concentrations in Italy, 2013–2015, using a
384 spatiotemporal land-use random-forest model. *Environment International* 124, 170–179.

385 Tai, A.P., Mickley, L.J., Jacob, D.J., 2010. Correlations between fine particulate matter (PM_{2.5}) and meteorological variables in the United States:
386 Implications for the sensitivity of PM_{2.5} to climate change. *Atmospheric Environment* 44, 3976–3984.

387 Tang, G., Zhang, J., Zhu, X., Song, T., Münkel, C., Hu, B., Schäfer, K., Liu, Z., Zhang, J., Wang, L., Xin, J., Suppan, P., Wang, Y., 2016. Mixing
388 layer height and its implications for air pollution over Beijing, China. *Atmospheric Chemistry and Physics* 16, 2459–2475.

389 Thompson, M.L., Reynolds, J., Cox, L.H., Guttorp, P., Sampson, P.D., 2001. A review of statistical methods for the meteorological adjustment of
390 tropospheric ozone. *Atmospheric Environment* 35, 617–630.

391 Vogelesang, D.H.P., Holtzlag, A.A.M., 1996. Evaluation and model impacts of alternative boundary-layer height formulations. *Boundary-Layer
392 Meteorology* 81, 245–269.

393 Vu, T.V., Shi, Z., Cheng, J., Zhang, Q., He, K., Wang, S., Harrison, R.M., 2019. Assessing the impact of clean air action on air quality trends in
394 Beijing using a machine learning technique. *Atmospheric Chemistry and Physics* 19, 11303–11314.

395 Xiang, Y., Zhang, T., Liu, J., Lv, L., Dong, Y., Chen, Z., 2019. Atmosphere boundary layer height and its effect on air pollutants in Beijing during
396 winter heavy pollution. *Atmospheric Research* 215, 305–316.

397 Xie, Y., Dai, H., Dong, H., Hanaoka, T., Masui, T., 2016. Economic impacts from PM_{2.5} pollution-related health effects in China: a provincial-level
398 analysis. *Environmental Science & Technology* 50, 4836–4843.

399 Zhang, Q., Ma, X., Tie, X., Huang, M., Zhao, C., 2009. Vertical distributions of aerosols under different weather conditions: Analysis of in-situ
400 aircraft measurements in Beijing, China. *Atmospheric Environment* 43, 5526–5535.

401 Zhang, S., Chen, S.X., Guo, B., Wang, H., Lin, W., 2020. Regional air-quality assessment that adjusts for meteorological confounding. *Scientia
402 Sinica Mathematica* 50, 527–558. The English translation is available from [https://songxichen.com/uploads/Files/Publication/
403 SCM-2019-0368R2_Proof_hi.pdf](https://songxichen.com/uploads/Files/Publication/SCM-2019-0368R2_Proof_hi.pdf).

404 Zhang, S., Guo, B., Dong, A., He, J., Xu, Z., Chen, S.X., 2017a. Cautionary tales on air-quality improvement in Beijing. *Proceedings of the Royal
405 Society A* 473, 20170457.

406 Zhang, W., Guo, J., Miao, Y., Liu, H., Song, Y., Fang, Z., He, J., Lou, M., Yan, Y., Li, Y., Zhai, P., 2018. On the summertime planetary boundary
407 layer with different thermodynamic stability in China: A radiosonde perspective. *Journal of Climate* 31, 1451–1465.

408 Zhang, X., Zhang, X., Chen, X., 2017b. Happiness in the air: How does a dirty sky affect mental health and subjective well-being? *Journal of
409 Environmental Economics and Management* 85, 81–94.