# Estimation of high-resolution PM2.5 over Indo-Gangetic Plain by fusion of satellite data, meteorology, and land use variables

Alaa Mhawish, Tirthankar Banerjee, Meytar Sorek-Hamer, Muhammad Bilal, Alexei Lyapustin, Robert B. Chatfield, and David Broday

**Just Accepted**

1    **Estimation of high-resolution PM$_{2.5}$ over Indo-Gangetic Plain by fusion of satellite**
2    **data, meteorology, and land use variables**

3    Alaa Mhawish[1,2,3], Tirthankar Banerjee[3,4]*, Meytar Sorek-Hamer[1,2], Muhammad Bilal[5], Alexei I. Lyapustin[6],
4    Robert Chatfield[2], David M Broday[7]

5    [1] Universities Space Research Association (USRA), CA, USA
6    [2] NASA Ames Research Center, Moffett Field, CA, USA
7    [3] Institute of Environment and Sustainable Development, Banaras Hindu University, Varanasi, India
8    [4] DST-Mahamana Centre of Excellence in Climate Change Research, Banaras Hindu University, Varanasi, India
9    [5] School of Marine Sciences, Nanjing University of Information Science and Technology, Nanjing, China
10   [6] NASA Goddard Space Flight Center, Greenbelt, MD, USA
11   [7] Civil and Environmental Engineering, Technion, Haifa, Israel
12
13   *Correspondence to: T. Banerjee (tb.iesd@bhu.ac.in; tirthankaronline@gmail.com)

14

15

16

17

18   **Key points:**


19   1. High-resolution MAIAC AOD-based PM$_{2.5}$ was estimated over the Indo-Gangetic Plain; fusing satellite
20   data, land-use variables & meteorology.

21   2. Random forest based AOD-PM model estimates were able to capture and quantify the PM$_{2.5}$ variability
22   at a sub-urban scale.

23   3. Comparatively high PM$_{2.5}$ concentrations were evident over central and lower IGP, mediated by land-
24   use and local meteorology.

25
26

27

28

29

30

31

32

33                                              **ABSTRACT**

34      Very high spatially resolved satellite-derived ground-level $PM_{2.5}$ concentrations have multiple potential

35      applications especially in air quality modelling, epidemiological and climatological research. Satellite-

36      derived aerosol optical epth (AOD), and columnar water vapor (CWV), meteorological parameters, and

37      land use data were used as variables within a linear mixed effect model (LME) and a random forest (RF)

38      model, to predict daily ground-level concentrations of $PM_{2.5}$ at 1km×1km grid across the Indo-Gangetic

39      Plain (IGP) in South Asia. The RF model exhibited superior performance and higher accuracy than the LME

40      model, with higher cross-validated explained variance ($R^2$=0.87) and lower relative prediction error

41      (RPE=24.5%). The RF model revealed improved performance metrics for increasing averaging periods,

42      from daily to weekly, monthly, seasonal, and annual means, which supports using it to estimate $PM_{2.5}$

43      exposure metrics across the IGP at varying temporal scales (i.e. both short and long terms). The RF-based

44      $PM_{2.5}$ estimates show high $PM_{2.5}$ levels over the middle and lower IGP, with the annual mean exceeding

45      110µg/m$^3$. Seasonally, winter was the most polluted season while monsoon was the cleanest. Spatially,

46      the middle and lower IGP showed poorer air quality compared to the upper IGP. In winter, the middle and

47      lower IGP experience very poor air quality, with mean $PM_{2.5}$ concentrations >170µg/m$^3$.

48      **Keywords**: Aerosols; Machine Learning; Random Forest; Mixed effect model; MAIAC; IGP.

49      **1.   Introduction**

50      Airborne fine particulate matter with an aerodynamic diameter less than 2.5 µm ($PM_{2.5}$) have been

51      associated with many adverse health effects, especially with cardiovascular and respiratory diseases[1].

52      Numerous epidemiological studies associate exposure to $PM_{2.5}$ with different health outcomes[2–7].

53      Recently, the World Health Organization estimated around 4.2 million deaths were attributable, globally,

54      to air pollution[8]. However, most of the epidemiological studies have been conducted in major urban areas,

55      where air quality monitoring is denser, rather than in small cities and rural areas. In South Asia, the air

56      quality monitoring stations are sparsely distributed and are found mainly in major cities, such as in Delhi,

57      Mumbai, Dhaka, and in state capitals. In the suburban and rural areas where a major fraction of the

58      population resides, and the $PM_{2.5}$ levels are as high as in urban areas[9–10]; there are only a few air quality

59      monitoring stations, and in some regions, there are none at all. Health risk assessments of $PM_{2.5}$ exposure

60      across highly populated and polluted areas of the Indo-Gangetic Plain (IGP) are severely constrained by

61      the sparse air quality monitoring stations and limited availability of particulate measurement data[11].

62    Understating the $PM_{2.5}$ spatial and temporal distribution therefore, is essential to improve understanding

63    of its impact on human health and regional climate.

64    Satellite remote sensing has the capability to provide high spatially resolved aerosol optical depth

65    measurements with daily global coverage, which can be used to predict ground-level particulate

66    concentration[12,13]. The aerosol optical depth (AOD) is a measure of the extinction of solar radiation by

67    aerosols in the atmospheric column, from the earth's surface to the top of the atmosphere. In contrast,

68    $PM_{2.5}$ is the mass concentration of fine particulate matter measured near the surface. Since both measures

69    (i.e. AOD, $PM_{2.5}$) are affected by the amount of suspended particles in the air, it is commonly assumed

70    that a correlation between AOD and $PM_{2.5}$ can be established and that AOD can be used to predict ground-

71    level $PM_{2.5}$ concentrations after accounting for factors that may interfere with the relationship, e.g. time-

72    varying parameters (RH, temperature, wind speed, etc.). Satellite-retrieved AOD has been widely used to

73    predict ground-level $PM_{2.5}$ concentrations, especially over the areas where ground monitoring stations

74    are not available. For example, in the last decade, satellite-retrieved AOD from different satellite-borne

75    sensors has been used for predicting ground-level $PM_{2.5}$ at varying spatial resolutions, including

76    instruments onboard Low Earth Orbit (LEO) satellites, such as Moderate Resolution Imaging

77    Spectroradiometer (MODIS)[14–19], Multiangle Imaging SpectroRadiometer (MISR)[9,20–22], and Visible Infrared

78    Imaging Radiometer Suite (VIIRS)[23–25], as well as instruments having a Geostationary Earth Orbit (GEO;

79    with only local coverage) satellites such as Himawari[26,27], and GOES[28,29]. In parallel, a new operational

80    MODIS aerosol retrieval algorithm named MultiAngle Implementation of Atmospheric Correction (MAIAC)

81    has been gaining attention as it provides AOD retrievals at very high spatial resolution (1km grid) with

82    global coverage and with each instrument possessing a daily revisit period. MAIAC AOD retrievals have

83    high capability in identifying fine aerosols, emission sources, and aerosol hotspots[30–32]. Hence, it has been

84    widely studied and found to be a powerful predictor of ground-level $PM_{2.5}$ concentrations compared to

85    other AOD products with coarser resolution[33–40]. However, several factors such as meteorology and

86    aerosol types can influence the relationship between AOD and $PM_{2.5}$[41]. Several studies suggested

87    considering other influential factors (such as meteorological variables, land use parameters, and aerosol

88    types) in the AOD-PM modeling to improve the $PM_{2.5}$ prediction using AOD measurements[42 35,39].

89    Various statistical models have been explored to establish the relation between satellite-retrieved

90    AOD and ground-level $PM_{2.5}$, extending from simple multivariate regression models[13] to more advanced

91    statistical models such as linear mixed effect models[14,15,37,43–45], geographically weighted regression

92    models[46–48], generalized additive models[22,49], and other nonlinear models[49–51]. Moreover, some studies

93    used multiple-stage models to address the spatiotemporal variations in the AOD-PM$_{2.5}$ relationship for

94    more accurate PM$_{2.5}$ predictions[34,39,52].

95        Recently, machine learning algorithms have been also applied to predict ground-level PM$_{2.5}$

96    concentrations[27,53–55]. Unlike traditional statistical models, machine learning algorithms have the ability to

97    use a large number of predictors with a few prior assumptions, thus enhancing their predictive power and

98    enabling to capture the complexity in the AOD-PM$_{2.5}$ relationships[56]. Ensemble models such as Random

99    Forest (RF) and the Gradient Boosting (GB) models combine weak learners (multiple models) to obtain

100   more accurate and robust models[57]. RF models have been successfully used to predict PM$_{2.5}$ over several

101   regions, such as in China[53], USA[55] and Italy[40]. In India, very few studies have been conducted to predict

102   ground-level PM$_{2.5}$ concentration using satellite AOD data. For example, Dey et al.[9] and Chowdhury et al.[58]

103   have been used AOD data obtained from MISR and MAIAC AOD respectively, to predict PM$_{2.5}$

104   concentration by multiplying the AOD with conversion factor obtained from GEOS-Chem chemical

105   transport model. Recently, Mandal et al.[59] implemented multiple-stage modeling including statistical

106   model and machine learning algorithm to predict PM$_{2.5}$ in the national capital of India using satellite data,

107   and land use, meteorological data, and population. However, to the best of our knowledge, no study that

108   reports PM$_{2.5}$ prediction has been conducted across South Asia using an advanced statistical or machine

109   learning model for PM$_{2.5}$ prediction on a regional scale.

110       In this study, both a statistical model (LME) and a machine learning algorithm (RF) were used to

111   predict, for the first time, high spatially resolved (1 km) ground-level PM$_{2.5}$ concentrations over the IGP

112   region, India; with the MAIAC AOD as an independent variable. The main objective of this study was to

113   examine how accurate can a machine learning model that uses the above satellite-based AOD product be

114   for estimating ground-level PM$_{2.5}$ concentrations. In response to this task, we first, compared the

115   performance of the RF model against that of an LME model, and the more accurate model was used for

116   PM$_{2.5}$ prediction. Next, we have studied the spatiotemporal variation of the estimated PM$_{2.5}$ across the

117   IGP region and identified regional hotspots. The dataset and model details used are described in section

118   2, the results are presented in section 3, and followed by a discussion in section 4.

119   **2.   Data and Method**
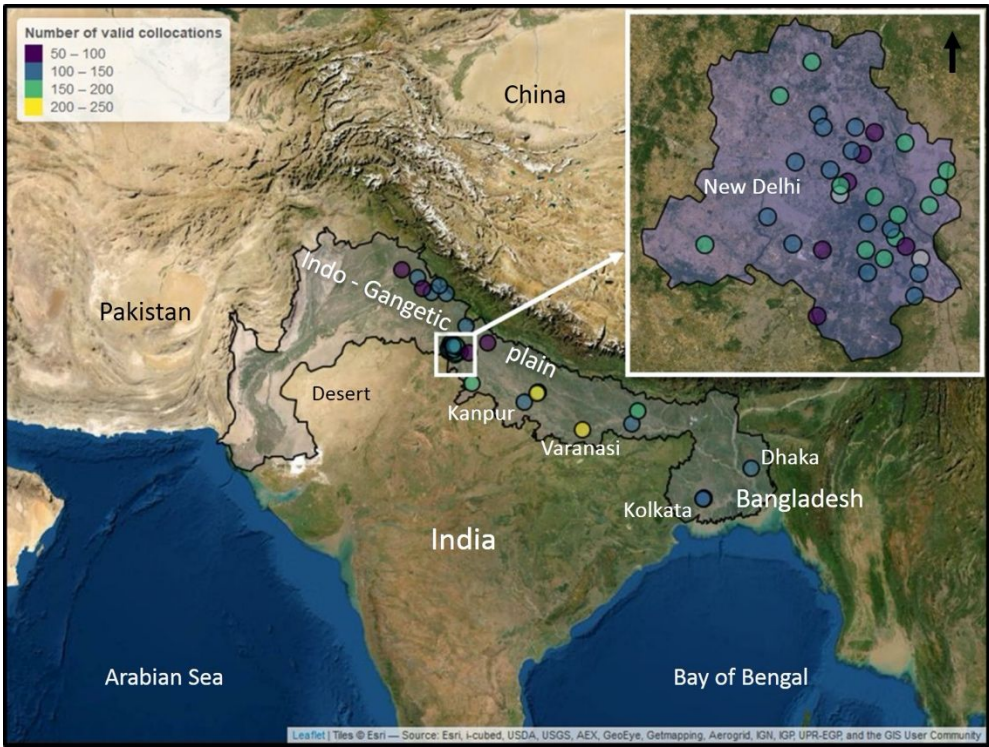
120   **2.1 Study region**

121

122    Figure 1: Map of the study region and the spatial distribution of PM$_{2.5}$ monitoring stations. The colors

123    represent the number of PM$_{2.5}$ and independent variables collocations. The shaded area represents the

124    IGP. The area within the box represents the monitoring stations in Delhi.

125        The study area covers the IGP, which stretches west from Pakistan across Northern India to the east

126    of the Bay of Bengal and Bangladesh (Fig. 1). The IGP region is densely populated and accommodates

127    nearly 13% (>800 million) of the world population. The rapid economic and population growth across the

128    region is associated with a wide range of anthropogenic activities, including biomass/-waste burning,

129    industries, and vehicular emissions, resulting in significant particulate matter pollution across the region.

130    The region is considered to be one of the aerosol hotspots and is characterized by a persistent high aerosol

131    loading throughout the year[60–62].

132    **2.2 Ground-based PM$_{2.5}$ Measurements**

133        Daily mean PM$_{2.5}$ concentrations were obtained from a total of 64 air quality monitoring stations

134    across the IGP from July 1$^{st}$, 2018, to June 30$^{th}$, 2019. Specifically, PM$_{2.5}$ data were obtained from 61

135    monitoring stations of the Central Pollution Control Board (CPCB) (https://app.cpcbccr.com/ccr/#/caaqm-

136    dashboard-all/caaqm-landing) and 3 PM$_{2.5}$ monitoring stations operated by the US Consulate in Delhi,

137    Kolkata, and Dhaka. The spatial distribution of the air quality monitoring stations across the region is

138    sparse and varying. For example, almost 50% of the monitoring stations are located in New Delhi, with

139    the rest distributed across the major cities. None of the air quality monitoring stations were in rural areas.

140    Measurement of $PM_{2.5}$ concentrations at all the monitoring stations is done with beta gauge

141    attenuation monitors (BAM-1020; Met One Instruments) that report hourly mean $PM_{2.5}$ concentrations.

142    We calculated the daily mean $PM_{2.5}$ concentrations after applying strict quality control procedures to

143    remove abnormal observation. Only days with more than 14 hourly measurements (60%) were used in

144    the analysis. The geographical location of the monitoring stations and the data availability of the valid

145    MAIAC AOD and $PM_{2.5}$ collocations are shown in Figure 1.

146    **2.3 MODIS MAIAC Products**

147    MAIAC is a relatively new operational MODIS-based aerosol retrieval algorithm that retrieves aerosol

148    properties and columnar water vapor at 1 km spatial resolution over land surface except for snow and

149    ice[32]. The MAIAC aerosol products have higher spatial resolution compared to other operational MODIS

150    aerosol products based on the Dark Target[63] (DT) and the Deep Blue[64] (DB) algorithms. Several validation

151    studies showed that the MAIAC algorithm improves aerosol retrieval accuracy, especially over bright

152    surfaces such as urban areas and dry land[30,32]. Several reasons make MAIAC significantly superior over

153    other operational MODIS algorithms: (a) the high spatial resolution (1km) compared to DT and DB (10km

154    and 3km) that allows to distinguish fine spatial features and to enhance spatial coverage[30], (b) high

155    retrieval accuracy over both dark and bright surfaces[30], and (c) MAIAC'S capability to retrieve AOD for

156    different aerosol types while discriminating among absorbing fine (smoke) and coarse (dust) aerosols[32].

157    In this study, the combined Terra and Aqua MAIAC product (MCD19A2;

158    https://ladsweb.modaps.eosdis.nasa.gov/) was used to extract Terra and Aqua AOD at 550nm with an

159    aerosol type (compositional) label (dust, smoke, and background), and CWV. Only the highest quality data,

160    designated with the Quality Assurance (QA) cloud mask value "clear", were used.

161    The spatial coverage of Terra (~10:30 am overpass time) and Aqua (~01:30 pm overpass time) MAIAC

162    AOD varies due to the diurnal cycle of cloud cover[30], meteorological conditions (mainly the lower

163    atmospheric boundary layer[65]) and the daily varying anthropogenic activities[66]. Therefore, a combined

164    MAIAC AOD product from both Terra and Aqua, can enhance the spatial and temporal coverage and

165    provide a more representative AOD that accounts for both the morning (Terra) and afternoon (Aqua) time

166    windows, from ~10:00 am until 02:00 pm local time. Nevertheless, if one of the two values (Aqua or Terra

167    AOD) is missing the combined AOD product will be biased towards either the morning or the afternoon

168    retrieval. To eliminate this bias, missing Terra AOD were predicted from Aqua AOD, and vice versa, by

169    fitting seasonal linear regression models to both the Aqua and Terra AOD. Table S1 shows the seasonal

170    regression equations and correlation coefficients of each season for both the AOD and CWV. The number

171    of the available combined AOD product increased by 22% and 24% compared to Terra or Aqua only AOD

172    retrievals, respectively. The $R^2$ of the seasonal regressions between Aqua and Terra MAIAC AOD ranged

173    from 0.63-0.79 ($p < 0.001$).

174    **2.4 Meteorological data**

175    Meteorological variables, including the ambient temperature at 2 m a.g.l. (temp; K), surface pressure

176    (SF; hPa), wind field at 10 m a.g.l. (Wind Speed (WS); m s$^{-1}$, and Wind Direction (WD); $^o$), relative humidity

177    (RH; %), and the planetary boundary layer height (PBLH; m), were obtained from the European Center for

178    Medium-Range    Forecast    (ECMWF)    atmospheric    reanalysis    ERA-Interim    products

179    (https://apps.ecmwf.int/datasets/data/interim-full-daily/levtype=sfc/). The spatial resolution of ERA-

180    Interim is 12.5km, and its temporal resolution is 6h except for the PBLH that is provided in every 3h. All

181    the meteorological variables were averaged over the time window corresponding to the Terra and Aqua

182    overpass times.

183    **2.5 Auxiliary data**

184    Level 3 Terra and Aqua MODIS 16-day composite Normalized Difference Vegetation Index (NDVI) data

185    (MxD13A2, x is O or Y for Terra and Aqua, respectively) at 1km spatial resolution were used in this study

186    as a proxy for the land use parameter. The NDVI data are reported every 16 days for both Terra and Aqua

187    but with 8 days of difference between them (Terra reports on day 001 while Aqua reports on day 009).

188    This corresponds to having a measurement in every 8 days. Elevation (Elev) data were obtained from the

189    Shuttle Radar Topography Mission (SRTM) database (http://srtm.csi.cgiar.org/srtmdata/) at 30m spatial

190    resolution[67] and used as spatial predictors.

191    **2.6 Data processing and integration**

192    For predicting $PM_{2.5}$ concentration at 1km spatial resolution using MAIAC AOD and other temporal

193    and spatial predictors, the spatial resolution of all the predictors should be consistent and matched with

194    the MAIAC AOD grid. Therefore, all the meteorological data and auxiliary data were re-projected and

195    gridded to match the MAIAC AOD fixed grid. In particular, the Terra and Aqua MODIS NDVI were gridded

196    to a 1km, and then the combined Terra and Aqua daily NDVI was calculated using the temporally

197    interpolated spline function technique. The meteorological data were gridded to a 1km grid using bi-linear

198    interpolation and temporal subsets that matched the Terra and Aqua overpass times (~10 am–2 pm local

199    time). Similarly, elevation data were also gridded to 1km. The total number of spatial and temporal PM$_{2.5}$,

200    satellite- and meteorological data was 8233 matched up collocations, distributed over 356 days from July

201    1$^{st}$, 2018 to June 30$^{th}$, 2019.

202    **2.7 Model Development**

203    **2.7.1 Linear Mixed Effect (LME) Model**

204         Linear mixed effect (LME) models have been widely used to estimate PM$_{2.5}$ concentrations based

205    on satellite-derived AOD[68] since it controls for the inherent day-to-day variability in the relationship

206    between AOD and PM$_{2.5}$. The AOD-PM$_{2.5}$ relationship is expected to be influenced by time-varying

207    parameters such as RH, PBLH, Temp, and the optical properties of the particles and their vertical

208    distribution[41]. Therefore, considering the daily variability in the AOD-PM$_{2.5}$ relationship is essential to

209    improve the correlation between the AOD and PM$_{2.5}$. Hence, allowing for a day-specific random slope and

210    intercept enables to examine the day-to-day variability in the AOD- PM$_{2.5}$ relationship.

211         In this study, we developed a nested day- and month-specific random effect based on all the days

212    with valid AOD-PM$_{2.5}$ collocations (the days with less than three collocations were removed from the

213    dataset). The LME model structure is expressed by the following equation (Eq. 1):

214    $$PM_{ij} = (\alpha_0 + (\alpha_{day} + \alpha_{month})) + (\beta_0 + (\beta_{day} + \beta_{month})) \times AOD_{ij} + \beta_1\,CWV_{ij} + \beta_2\,WS_{ij}$$
215    $$+\ \beta_3\,RH_{ij} + \beta_4\,PBLH_{ij} + \beta_5\,SP_{ij} + \beta_6\,WD_{ij} + \beta_7\,NDVI_{ij} + \beta_8\,Temp_{ij} + \beta_9\,Elev_{ij} + \varepsilon_{ij}$$

216                                                                                                      Eq. (1)

217    where $PM_{ij}$ and AOD$_{ij}$ are the PM$_{2.5}$ concentration and MAIAC AOD at monitoring site $i$ on day $j$; $\alpha_0$ and $\beta_0$

218    are the fixed intercept and slope, respectively; $\alpha_{day}$, $\alpha_{Month}$, $\beta_{day}$ and $\beta_{month}$ are the day- and month-

219    specific random intercept and slope, respectively; CWV, WS, WD, RH, PBLH, SP, NDVI, Temp, and Elev are

220    the corresponding auxiliary variables at grid $i$ and day $j$ (and their corresponding fixed slopes); and $\varepsilon_{ij}$ is

221    the error term at site $i$ and on day $j$.

222         The day- and month-specific intercepts and slopes allow the model to control day-to-day and

223    monthly variability in the relationship between PM$_{2.5}$ and AOD. Spatial predictors such as NDVI, and

224    Elevation were found to be significantly correlated with the PM$_{2.5}$ and, therefore, were included in the

225    model. All the variables were tested and only the significant ones were used during the model fitting.

226    **2.7.2    Random Forest (RF) model**

227     Random forest is an ensemble learning (algorithm) that aggregates a large number of decision trees

228     which were created independently using the bootstrap resampling method[69]. The bagging (bootstrap

229     aggregation) technique allows to reduce the variance of the estimated prediction by averaging the

230     regression results from all decision trees. Each node of the tree splits into two daughter nodes using the

231     best split from the randomly selected variables[69]. Aggregating weak learners into a strong learner leads

232     to a final model of enhanced performance. Moreover, the random forest model provides an estimate for

233     the importance of each variable by measuring the increase in the prediction error (decrease in the

234     accuracy score) of the final model after performing variable permutations. Here, the mean decrease

235     accuracy is calculated by the permutation scheme of Breiman[70]. In the random forest algorithm, the main

236     two variables that have the major effect at each level (bifurcation) on the model accuracy are mostly the

237     ones used to split the residual subset at each node (mtry) and to select the number of trees in the forest

238     (ntree). In this experiment, we found that the best model accuracy was obtained for mtry=12 and

239     ntree=1500. The unscaled variables importance[71] of the final model is reported in Table S2.

240     **2.8 Evaluation of Models**

241     To evaluate the performance of the developed models across the IGP, we adopted two 10-fold cross-

242     validations (CV) approaches a site-based CV, and a sample-based CV. The 10-fold CV method[72] randomly

243     split the database into ten subsets, each containing 10% of the data. In each round, the model trains on

244     nine subsets (90% of the data) and predicts the $10^{th}$ subset, with the predictions evaluated against the

245     true data. The process is repeated ten times thus ensuring that every subset has been evaluated. In the

246     site-based CV, the database is split according to the monitoring sites into ten subsets, each containing

247     ~10% of the data. As such, each subset contains different monitoring stations. In each round, one subset

248     is held-out and PM levels at the sites it contains are predicted using the model that has been developed

249     based on data from the other sites. The model is evaluated by comparing its predictions in the held-out

250     sites against the true observations (which have not been used for the model development). This process

251     is repeated with each subset held-out in turn. The site-based CV is used for assessing how well the model

252     performs over regions that do not have monitoring sites, such that the prediction must be done by

253     applying a model that has been developed (and evaluated) over another region. In the sample-based CV,

254     the same procedure is performed using the whole database without accounting for the monitoring site it

255     comes from. As such, the sample-based CV is used for a general assessment of the model performance

256     for filling data gaps in both space and time in regions were monitoring sites do exist. The performance of

257     the CV predictions has been examined using several statistical metrics, including the Root Mean Squared

258    Error (RMSE), Relative Prediction Error (RPE; Eq. 2), coefficient of determination ($R^2$), Mean Prediction

259    Error (MPE; Eq. 3), and the slope (*b*) and intercept (*a*) of the linear regression between the predicted and

260    observed $PM_{2.5}$. The RPE and the MPE are calculated as:

261    $$\text{RPE} = \frac{\text{RMSE}}{\overline{\text{PM}_{2.5}}} \times 100$$                    Eq. (2)

262    $$\text{MPE} = \frac{1}{N} \Sigma_{i=1}^{N} |\text{Predicted}_i - \text{Observed}_i|$$                    Eq. (3)

263

264    **3.  Results and Discussion**

265    **3.1 Descriptive statistics**

266        The histograms and descriptive statistics for all the dependent and independent variables used for

267    the model development are illustrated in Figure S1 and Tables S3. The annual mean $PM_{2.5}$ over the entire

268    region was 114.49 ± 76.65 μg/m³ (*N*=8,233), and the seasonal mean $PM_{2.5}$ were winter: 170.16 ± 88.46

269    μg/m³, postmonsoon: 150.69 ± 73.16 μg/m³, premonsoon: 77.59 ± 34.38 μg/m³, and monsoon: 58.00 ±

270    21.98 μg/m³. The overall mean AOD was 0.57 ± 0.39, and the seasonal means were winter: 3.28 ± 0.40,

271    postmonsoon: 0.77 ± 0.55, premonsoon: 0.40 ± 0.20 and monsoon: 3.13 ± 0.28 (Table S3). Notably, while

272    the highest $PM_{2.5}$ was observed in the winter and the lowest $PM_{2.5}$ was observed in the monsoon, the AOD

273    showed much smaller variation with the highest retrievals during the postmonsoon and the lowest

274    retrievals during the pre-monsoon seasons. The seasonal discrepancies between AOD and $PM_{2.5}$, in

275    particular, the low $PM_{2.5}$ concentrations but high AOD values during the monsoon season, are attributed

276    to the abundance in water vapor in the atmospheric column during monsoon (CWV =3.32 ± 0.65), which

277    favors hygroscopic growth of the aerosol particles[30,61]. Hygroscopic growth of aerosol particles enhances

278    scattering, thus resulting in higher AOD[73]. In contrast, $PM_{2.5}$ is measured near the surface at a fixed RH of

279    <40% and does not reflect hygroscopic growth as in the free air. Both the AOD and $PM_{2.5}$ data showed a

280    similar unimodal distribution, with the correlation coefficient between the daily mean $PM_{2.5}$ and the

281    combined Aqua and Terra MAIAC AOD being *r* =0.47 (*p* < 0.0001). The variables used in this study, i.e.

282    meteorological variables, boundary layer height, and land use and the land cover attributes, were all

283    significantly correlated (*p* < 0.0001) with the $PM_{2.5}$ (Table S4).

284        Furthermore, the variance inflation factors (VIF) was used to quantify the collinearity among the

285    predictors, which could affect the model performance. A VIF value of 10 was set as the threshold for

286    collinearity. All the VIF values were <10, i.e. showing little to nil collinearity (Table S5).

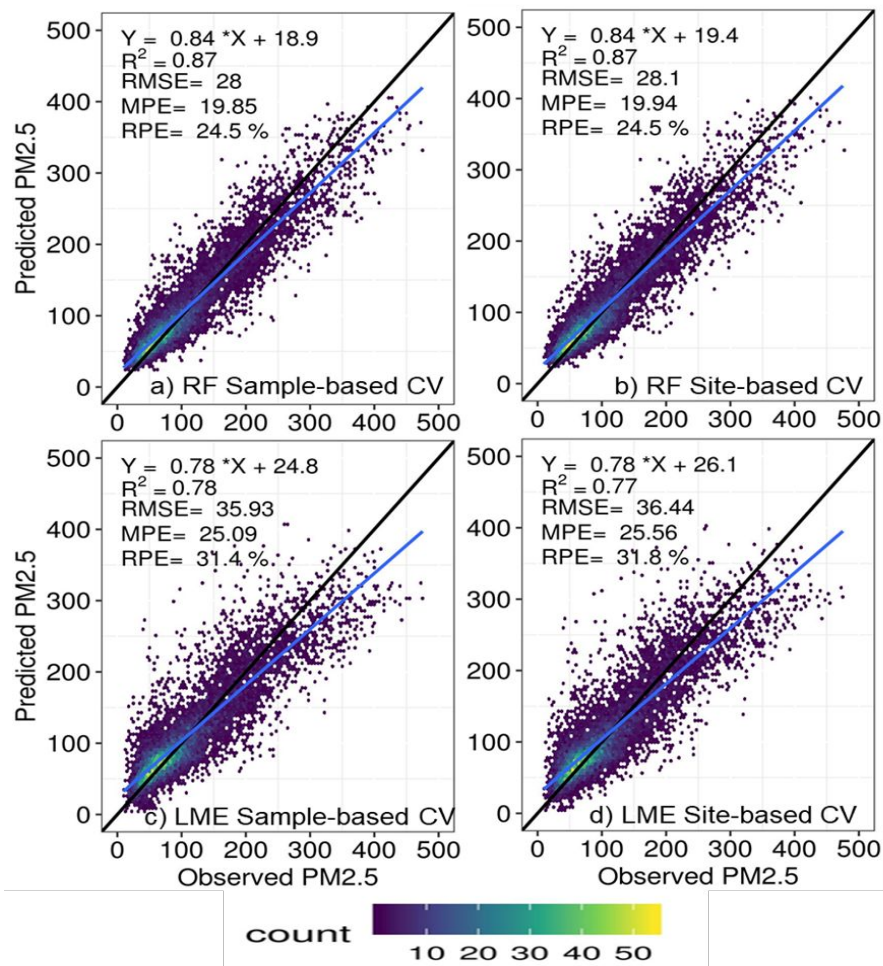287    **3.2  Models fitting and evaluation**



288

289    Figure 2: Scatterplot of the cross-validation results. Left column: sample-based cross-validation, right

290    column: site-based cross-validation, upper row: RF model, and lower row: LME model.

291        Figure 2 shows scatter plots of the sample-based and the site-based CV predicted vs observed daily

292    mean $PM_{2.5}$ for the LME and the RF models. Clearly, the RF model performed better, with $R^2$ of 0.87 and

293    RMSE of 28 μg/m³ (irrespective of the CV method applied), compared to the LME model ($R^2$ ~78%, RMSE

294    ~36 μg/m³). Both the RF and LME models tend to underestimate the ground-level $PM_{2.5}$ concentrations,

295    especially on highly polluted days ($PM_{2.5}$ >100 μg/m³), with the underestimation more severe when using

296    the LME model compared to the RF model. Since the sample-based and site-based CV methods resulted

297    in almost identical results, both the LME and RF models were apparently not over-fitted to the data,

298    suggesting a good spatial predictive power. Still, the RF model outperformed the LME in terms of accuracy,

299    having lower RPE (RF: 24.5%, LME: ~ 31.6%).

300        To evaluate the performance of the RF and LME models at different temporal averaging scales,

301    the weekly, monthly, seasonal, and annual $PM_{2.5}$ means were calculated based on daily predictions from

302    days for which >20% of the site-specific daily $PM_{2.5}$ predictions were available (figure S2 and Table S6).

303    Like on the daily scale, the RF model was more accurate than the LME model also on the weekly, monthly,

304    seasonal, and annual scales, with high $R^2$ (0.91-0.92), a slope close to unity (0.88-0.9), and a lower RPE
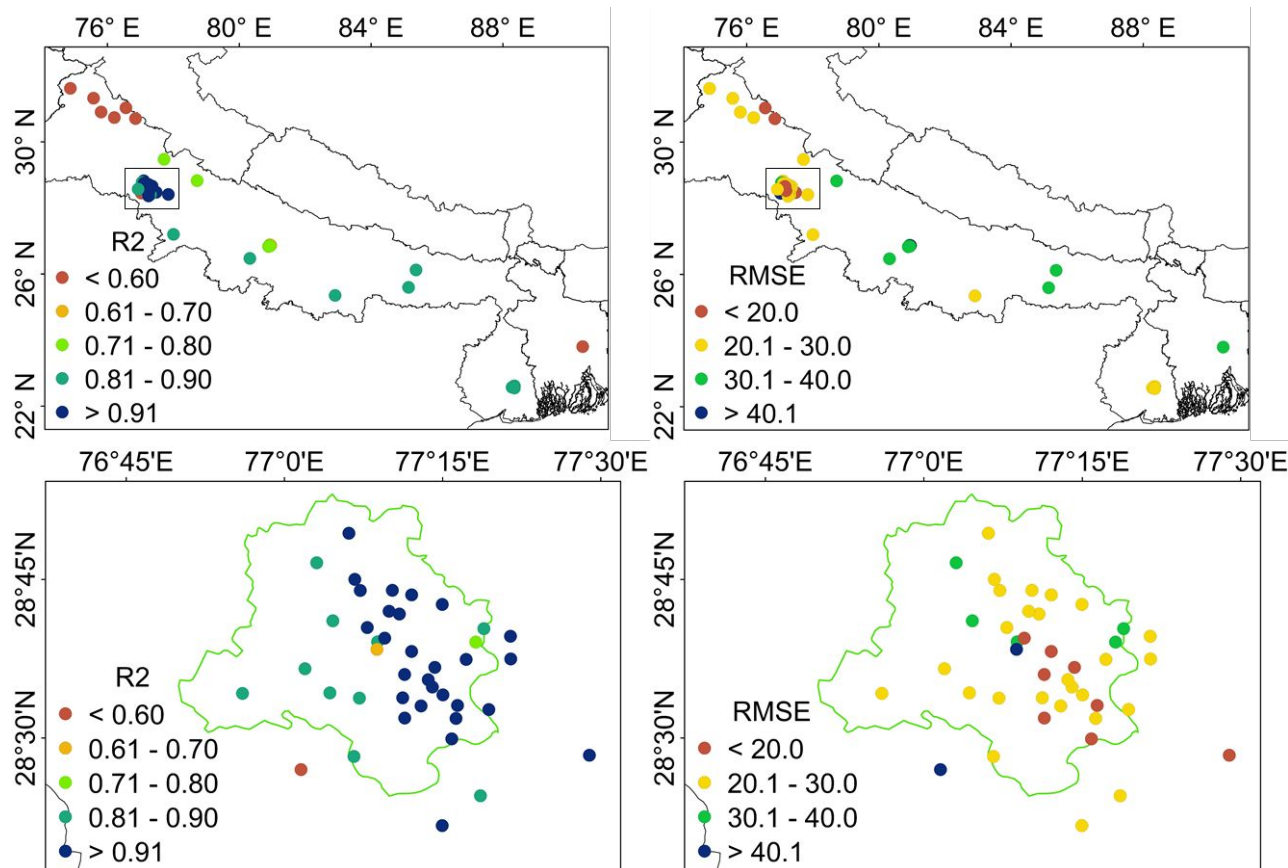
305    (monthly: 15.1%, seasonal: 13.9%, annual: 8.8%).



306

307    Figure 3: Spatial distribution of the $R^2$ (left panels) and RMSE (right panels) between the observed and

308    predicted $PM_{2.5}$ using the RF model. Upper row: across the whole IGP region, lower row: the Delhi area.

309    Figure 3 shows the spatial distribution of the site-based CV performance metrics of the RF model, with

310    Fig. 3(a, b) depicting the $R^2$ and RMSE across the IGP, respectively, and Fig. 3(c, d) focusing on the Delhi

311    city. The overall IGP mean $R^2$ was 0.81, with 85% of the stations showing $R^2$ >0.7. The lower $R^2$ (<0.6) were

312    found at the northwest IGP, which may be attributed to the small number of data points (collocations) in

313    this region due to limited $PM_{2.5}$ (only a few months), and may not represent the entire year. Similarly,

314    while the IGP average RMSE was 26.7 μg/m$^3$ a relatively high RMSE values (>30 μg/m$^3$) were evident in

315    few stations in the middle and lower IGP, attributed to the high PM$_{2.5}$ levels throughout the year (annual

316    average >150 μg/m$^3$). Since the model was trained with around 50% of the data obtained from stations in

317    Delhi, therefore the performance metrics were relatively better than the stations located in upper and

318    lower IGP. Overall, the RF model achieved satisfactory performance and was able to capture most of the

319    variability in the PM$_{2.5}$ across the region, with R$^2$ >0.7 in most of the monitoring stations.

320         The seasonal variation in aerosol sources and meteorological variables also affected the AOD-PM

321    model performance seasonally. In premonsoon and monsoon seasons, the IGP is affected by aerosols

322    transported by the southwest monsoon and frequently associated with higher wind speed and deeper

323    boundary layer. While in winter, the PM$_{2.5}$ primarily concentrates near the surface due to shallow

324    boundary layer and slower wind speed[61,62]. Model performances both in cold seasons i.e., winter (RPE:

325    20.9%) and postmonsoon (RPE: 22.3%) was also compared with warm seasons including premonsoon

326    (RPE: 28%) and monsoon (RPE: 26.5%) and shown in Table S7.  The larger slope of the fitting line for colder

327    seasons reflect a higher PM$_{2.5}$ that was concentrated near the surface due to a shallow PBL (788-1113m)

328    compared to the warmer seasons when the PBLH was relatively higher (1673-1901m)[74].

**3.3 Predicted PM$_{2.5}$ over IGP**

330    Figure 4 shows the annual mean satellite-based PM$_{2.5}$ estimates for IGP at 1 km grid resolution, as

331    derived from the RF model. The overall estimated annual mean PM$_{2.5}$ (July 1$^{st}$, 2018 to June 30$^{th}$, 2019)

332    was 112.7 μg/m$^3$, which exceeds the 40 μg/m$^3$ Indian National Ambient Air Quality Standards (NAAQS). In

333    particular, the middle and lower IGP regions experience higher PM$_{2.5}$ concentrations (>110 μg/m$^3$), with

334    around 79.3% of the area experiencing an annual mean PM$_{2.5}$ concentration between 110-150 μg/m$^3$. The

335    highest annual mean PM$_{2.5}$ was found over the state of Bihar, West of Bengal, and Bangladesh, with PM$_{2.5}$

336    concentrations exceeding 130 μg/m$^3$. The high PM$_{2.5}$ levels in the middle and lower IGP are most likely

337    due to the combined contributions of local sources and long-range transport from the upper IGP[61].
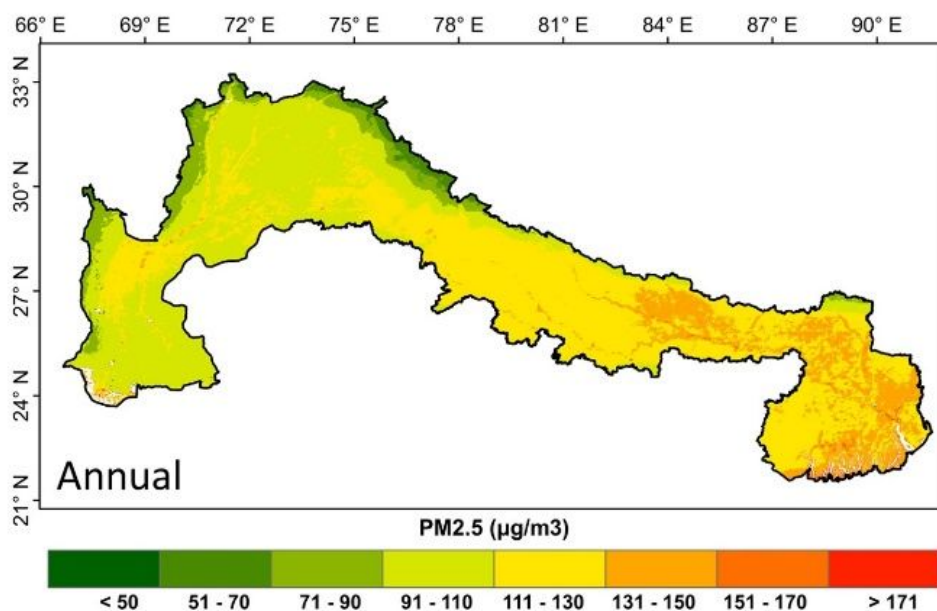
Figure 4: Spatial distribution of annual mean PM$_{2.5}$ estimates at 1 km grid resolution over IGP.

Seasonally, a significant variation is noted across the region with the highest PM$_{2.5}$ levels recorded in winter (DJF) (154 ± 22.4 μg/m$^3$) (Fig. S3). Spatial differences are also evident, with about 66% of the IGP exposed to PM$_{2.5}$ concentrations >150 μg/m$^3$ in the winter, and 25% of the IGP region (mainly the middle and lower IGP) experiencing PM$_{2.5}$ >170 μg/m$^3$. High PM$_{2.5}$ levels are also estimated during postmonsoon (ON; 128.8 ± 16.0 μg/m$^3$), with about 87% of the IGP exposed to PM$_{2.5}$ in the range 110 - 150 μg/m$^3$. The lowest PM$_{2.5}$ levels were estimated in the monsoon and premonsoon seasons, with mean PM$_{2.5}$, mean of 59.9 ± 5.9 μg/m$^3$ and 80.9 ± 9.5 μg/m$^3$, respectively.

Taking Delhi as an example for one of the most heavily PM$_{2.5}$-polluted metropolitans/megapolises in South Asia and the world, we also examined the capability of our model to capture PM$_{2.5}$ variability at the urban scale. The true color image (Fig. 5a) and the annual mean PM$_{2.5}$ estimates over Delhi (Figure 5b) show high PM$_{2.5}$ concentration in central and eastern Delhi – the most densely populated areas, and lower levels in southern Delhi; which is greener and not as densely populated, with an overall annual mean of PM$_{2.5}$ of 121.8 ± 7.4 μg/m$^3$, i.e. 8% higher than the IGP mean PM$_{2.5}$. These results suggest that both local particulate sources combined with long-range transport of aerosol from the north-west IGP, especially during stubble burning period [61, 75,76], could be captured by the model, which accounts for enhanced PM$_{2.5}$ concentrations in Delhi.

To critically examine a severe PM$_{2.5}$ condition, we selected the stubble burning episode, which occurs every year in November in the Punjab and Haryana states, and affects the whole northern India[76]. Figure 5(c, d) shows the spatial distribution of active fires on November 8th, 2018, obtained from the VIIRS

359    and MODIS (Aqua and Terra) sensors (https://firms.modaps.eosdis.nasa.gov/) together with the

360    estimated $PM_{2.5}$. Clearly, areas located downwind of the fire spots experienced higher $PM_{2.5}$ and the

361    model shows good sensitivity for capturing these high $PM_{2.5}$ areas, demonstrating the excellent capability

362    of the model to identify pollution sources in both space and time. Examining the model performance at

363    different Indian air quality categories, for $PM_{2.5}$ >60 µg/m$^3$ (moderate air quality) the CV $R^2$ was 0.84 and

364    the RPE was 21.8%, for $PM_{2.5}$ >90 µg/m$^3$ the $R^2$ was 0.79 and the RPE was 19.8%, and for $PM_{2.5}$ >120 µg/m$^3$

365    (very poor air quality) the $R^2$ and RPE were 0.71 and 20.82%, respectively (Fig. S4). The model performance

366    at low $PM_{2.5}$ (cleaner conditions) was poorer than at the polluted conditions (Fig. S4).



367
368    Figure 5: (a) RGB image, b) annual mean $PM_{2.5}$ over Delhi, c) Active fire counts in the northwest IGP

369    obtained from VIIRS and MODIS (Aqua and Terra) sensors on November 8$^{th}$, 2018, and d) $PM_{2.5}$ estimates

370    superimposed by the wind direction during the same day.

371    **4.  Discussion**

372         RF and LME models were developed to predict daily ground-level $PM_{2.5}$ concentrations across the IGP,

373    South Asia. A few studies have estimated $PM_{2.5}$ in the IGP region based on satellite retrieved AOD. Dey et

374    al.[9] used AOD retrievals from the MISR sensor to estimate ground-level $PM_{2.5}$ at a spatial resolution of 0.5°

375    x 0.5˚, using a scale factor obtained from the GEOS-Chem chemical transport model. Similarly, using a

376    scale factor from GEOM-Chem, Chowdhury et al.[58] estimated ground-level PM$_{2.5}$ concentrations during

377    the dry season (October–June) at a spatial resolution of 1km grid over the Delhi National Capital Region

378    (NCR). Similarly, over Delhi, Mandal et al.[59] estimated the PM$_{2.5}$ from 2010 to 2016 at 1km spatial grid

379    using the multi-stage prediction model. To the best of our knowledge, the current study is the first

380    regional-scale study in South Asia to predict daily PM$_{2.5}$ at a high spatial resolution (1km x 1km), using

381    satellite retrievals of AOD and CWV, together with meteorological and land use information, and applying

382    both the random forest (machine learning) algorithm as well as an advanced statistical model (LME). To

383    date, only few attempts were made to use satellite-based AOD for estimating ground-level PM$_{2.5}$ across

384    South Asia, as the region is severely constrained by the availability of quality surface monitoring data

385    which are essential for model calibration and validation. Taking advantage of the recently established air

386    quality monitoring network, we developed regional-scale models for estimating daily PM$_{2.5}$

387    concentrations over the IGP. The RF model exhibited adequate performance and higher accuracy than the

388    LME model, with better cross-validated explained variance (R$^2$ =0.87) and lower prediction error (RPE

389    =24.5%). Our RF model performed similarly to- or better than previous RF models that were developed

390    for China (R$^2$= 0.83-0.85, RPE=30.7%-35.9%)[53,54] and the USA (R$^2$=0.80, RPE=29.2%)[55]. The model showed

391    satisfactory predictive capability across the region with comparable site-based CV and sample-based CV

392    results. Moreover, the RF model also revealed high accuracy in estimating weekly (R$^2$ =0.91, RPE =17.7%),

393    monthly (R$^2$= 0.92, RPE= 15.5%), seasonal (R$^2$ =0.92, RPE =13.91%), and annual (R$^2$ =0.90, RPE =8.8%) mean

394    PM$_{2.5}$ levels. The high spatial resolution and low-bias of the PM$_{2.5}$ estimates (both weekly and monthly

395    mean) support using it in different research domains, especially in environmental epidemiology and

396    climatological studies.

397        Due to the lack of historical PM$_{2.5}$ records across the IGP, the year-to-year variability in PM$_{2.5}$

398    concentrations cannot be assessed. Similarly, the insufficient number of observations and the low PM$_{2.5}$

399    concentrations in the northwestern IGP resulted in poor model performance compared to other parts of

400    IGP. Indeed, the model results in greater accuracy when high PM$_{2.5}$ concentrations were experienced.

401        The modeled PM$_{2.5}$ map showed significant spatial and temporal variation across the IGP. Seasonally,

402    winter and postmonsoon are the more polluted seasons while the wet monsoon season is the cleaner

403    one. Anthropogenic activities such as open burning stubble during postmonsoon, and burning of biomass

404    and coal for heating and cooking, combined with shallow atmospheric boundary layer height, lead to

405    enhanced PM$_{2.5}$ concentrations during the winter and post-monsoon. The lower PM$_{2.5}$ concentrations in

406     the monsoon period are due to wet deposition, strong convection, and higher boundary layer heights.

407     Nonetheless, the $PM_{2.5}$ concentrations in all seasons were higher than the Indian NAAQS (annual average:

408     40 $\mu g/m^3$).

409     Spatially, the middle and lower IGP showed poor air quality compared to the upper IGP. In winter, the

410     middle and lower IGP experience very poor air quality, with mean $PM_{2.5}$ concentrations >170 $\mu g/m^3$. The

411     highly spatially resolved $PM_{2.5}$ estimates were found to have potential to identify $PM_{2.5}$ hotspots and to

412     study $PM_{2.5}$ on small scales, especially in urban areas. Our model performed well at the urban scale,

413     showing the good capability to capture spatial $PM_{2.5}$ variability.

414     Finally, the random forest machine learning algorithm showed high skill in predicting $PM_{2.5}$ by fusing

415     satellite aerosol products, meteorological models' output, and land use data. Future improvements of the

416     model may involve using richer land use parameters (i.e. the road network, vehicle volumes) and

417     emissions data (agricultural residues burning, industries emissions inventory, municipal solid waste

418     burning, etc.) which may be helpful to further improve the reliability of the AOD-PM model across the

419     Indo-Gangetic plain.

428     **Data availability**

429     MODIS MAIAC data is available at https://ladsweb.modaps.eosdis.nasa.gov. Modis Fire products are
430     obtained      from      Fire      Information      for      Resource      Management      System      (FIRMS)
431     (https://firms.modaps.eosdis.nasa.gov). All datasets were last accessed in November 2019.

432     **Author Contributions**

433     **AM**: methodology, formal analysis, review and writing draft manuscript; **MSH, MB, AL, RC, DB**:

434     methodology and interpretation; review and editing draft; **TB**: methodology and interpretation,

435     resources, review, writing and editing draft manuscript.

436    **Competing interests.** Authors declare that they have no conflict of interest.

437    **Supporting Information.** The supporting figures and tables are included in supplementary file.
438

439

440    6.    **Bibliography**

441    (1)    GBD 2013 Risk Factors Collaborators. Global, regional, and national comparative risk assessment
442          of 79 behavioural, environmental and occupational, and metabolic risks or clusters of risks in 188
443          countries, 1990-2013: a systematic analysis for the Global Burden of Disease Study 2013. *Lancet*
444          **2015**, *386*, 2287–2323.

445    (2)    Hayes, R.B., Lim, C., Zhang, Y., Cromar, K., Shao, Y., Reynolds, H.R., Silverman, D.T., Jones, R.R.,
446          Park, Y., Jerrett, M. and Ahn, J., PM2.5 air pollution and cause-specific cardiovascular disease
447          mortality. Int. J. Epidemiol. 2020, 49, 25–35.

448    (3)    Chowdhury, S.; Dey, S.; Smith, K. R. Ambient PM2.5 exposure and expected premature mortality
449          to 2100 in India under climate change scenarios. *Nat. Commun.* **2018**, *9*, 318.

450    (4)    Brook, R. D.; Newby, D. E.; Rajagopalan, S. Air pollution and cardiometabolic disease: an update
451          and call for clinical trials. *Am. J. Hypertens.* **2017**, *31*, 1–10.

452    (5)    Weber, S. A.; Insaf, T. Z.; Hall, E. S.; Talbot, T. O.; Huff, A. K. Assessing the impact of fine
453          particulate matter (PM2.5) on respiratory-cardiovascular chronic diseases in the New York City
454          Metropolitan area using Hierarchical Bayesian Model estimates. *Environ. Res.* **2016**, *151*, 399–
455          409.

456    (6)    Thurston, G. D.; Ahn, J.; Cromar, K. R.; Shao, Y.; Reynolds, H. R.; Jerrett, M.; Lim, C. C.; Shanley,
457          R.; Park, Y.; Hayes, R. B. Ambient Particulate Matter Air Pollution Exposure and Mortality in the
458          NIH-AARP Diet and Health Cohort. *Environ. Health Perspect.* **2016**, *124*, 484–490.

459    (7)    Lelieveld, J.; Evans, J. S.; Fnais, M.; Giannadaki, D.; Pozzer, A. The contribution of outdoor air
460          pollution sources to premature mortality on a global scale. *Nature* **2015**, *525*, 367–371.

461    (8)    World Health Organization. *Global action plan on physical activity 2018–2030: more active*
462          *people for a healthier world. Geneva:World Health Organization; 2018. Licence: CC BY-NC-SA 3.0*
463          *IGO.* ; 2018.

464    (9)    Dey, S.; Di Girolamo, L.; van Donkelaar, A.; Tripathi, S. N.; Gupta, T.; Mohan, M. Variability of
465          outdoor fine particulate (PM2.5) concentration in the Indian Subcontinent: A remote sensing
466          approach. *Remote Sens. Environ.* **2012**, *127*, 153–161.

467    (10)   Bilal, M., Nichol, J.E., Nazeer, M., Shi, Y., Wang, L., Kumar, K.R., Ho, H.C., Mazhar, U., Bleiweiss,
468          M.P., Qiu, Z. and Khedher, K.M. Characteristics of Fine Particulate Matter (PM2.5) over Urban,
469          Suburban, and Rural Areas of Hong Kong. Atmosphere 2019, 10, 496..

470  (11)  Murari, V.; Singh, N.; Ranjan, R.; Singh, R. S.; Banerjee, T. Source apportionment and health risk
471        assessment of airborne particulates over central Indo-Gangetic Plain. *Chemosphere* **2020**,
472        127145.

473  (12)  Wang, J. Intercomparison between satellite-derived aerosol optical thickness and $PM_{2.5}$ mass:
474        Implications for air quality studies. *Geophys. Res. Lett.* **2003**, *30*, 2095.

475  (13)  Gupta, P.; Christopher, S. A. Particulate matter air quality assessment using integrated surface,
476        satellite, and meteorological products: 2. A neural network approach. *J. Geophys. Res.* **2009**, *114*.

477  (14)  Xiao, Q.; Wang, Y.; Chang, H. H.; Meng, X.; Geng, G.; Lyapustin, A.; Liu, Y. Full-coverage high-
478        resolution daily PM2.5 estimation using MAIAC AOD in the Yangtze River Delta of China. *Remote*
479        *Sens. Environ.* **2017**, *199*, 437–446.

480  (15)  Sorek-Hamer, M.; Kloog, I.; Koutrakis, P.; Strawa, A. W.; Chatfield, R.; Cohen, A.; Ridgway, W. L.;
481        Broday, D. M. Assessment of PM 2.5 concentrations over bright surfaces using MODIS satellite
482        observations. *Remote Sens. Environ.* **2015**, *163*, 180–185.

483  (16)  Chudnovsky, A. A.; Koutrakis, P.; Kloog, I.; Melly, S.; Nordio, F.; Lyapustin, A.; Wang, Y.; Schwartz,
484        J. Fine particulate matter predictions using high resolution Aerosol Optical Depth (AOD)
485        retrievals. *Atmos. Environ.* **2014**, *89*, 189–198.

486  (17)  Zou, B.; Pu, Q.; Bilal, M.; Weng, Q.; Zhai, L.; Nichol, J. E. High-Resolution Satellite Mapping of Fine
487        Particulates Based on Geographically Weighted Regression. *IEEE Geosci. Remote Sensing Lett.*
488        **2016**, *13*, 495–499.

489  (18)  Bilal, M.; Nichol, J. E.; Spak, S. N. A new approach for estimation of fine particulate
490        concentrations using satellite aerosol optical depth and binning of meteorological variables.
491        *Aerosol Air Qual. Res.* **2017**, *17*, 356–367.

492  (19)  Chu, H.-J.; Bilal, M. PM2.5 mapping using integrated geographically temporally weighted
493        regression (GTWR) and random sample consensus (RANSAC) models. *Environ. Sci. Pollut. Res. Int.*
494        **2019**, *26*, 1902–1910.

495  (20)  Franklin, M.; Kalashnikova, O. V.; Garay, M. J. Size-resolved particulate matter concentrations
496        derived from 4.4 km-resolution size-fractionated Multi-angle Imaging SpectroRadiometer (MISR)
497        aerosol optical depth over Southern California. *Remote Sens. Environ.* **2017**, *196*, 312–323.

498  (21)  Franklin, M.; Chau, K.; Kalashnikova, O.; Garay, M.; Enebish, T.; Sorek-Hamer, M. Using Multi-
499        Angle Imaging SpectroRadiometer Aerosol Mixture Properties for Air Quality Assessment in
500        Mongolia. *Remote Sens (Basel)* **2018**, *10*, 1317.

501  (22)  Meng, X.; Garay, M. J.; Diner, D. J.; Kalashnikova, O. V.; Xu, J.; Liu, Y. Estimating PM 2.5''
502        speciation concentrations using prototype 4.4 km-resolution MISR aerosol properties over
503        Southern California. *Atmos. Environ.* **2018**, *181*, 70–81.

504  (23)  Wu, J.; Yao, F.; Li, W.; Si, M. VIIRS-based remote sensing estimation of ground-level PM 2.5
505        concentrations in Beijing–Tianjin–Hebei: A spatiotemporal statistical model. *Remote Sens.*
506        *Environ.* **2016**, *184*, 316–328.

507  (24)  Yao, F.; Si, M.; Li, W.; Wu, J. A multidimensional comparison between MODIS and VIIRS AOD in
508        estimating ground-level PM2.5 concentrations over a heavily polluted region in China. *Sci. Total*
509        *Environ.* **2018**, *618*, 819–828.

510  (25)  Yao, F.; Wu, J.; Li, W.; Peng, J. Estimating Daily PM2.5 Concentrations in Beijing Using 750-M
511        VIIRS IP AOD Retrievals and a Nested Spatiotemporal Statistical Model. *Remote Sens (Basel)*
512        **2019**, *11*, 841.

513  (26)  Wang, W.; Mao, F.; Du, L.; Pan, Z.; Gong, W.; Fang, S. Deriving Hourly PM2.5 Concentrations from
514        Himawari-8 AODs over Beijing–Tianjin–Hebei in China. *Remote Sens (Basel)* **2017**, *9*, 858.

515  (27)  Zang, L.; Mao, F.; Guo, J.; Wang, W.; Pan, Z.; Shen, H.; Zhu, B.; Wang, Z. Estimation of
516        spatiotemporal PM1.0 distributions in China by combining PM2.5 observations with satellite
517        aerosol optical depth. *Sci. Total Environ.* **2019**, *658*, 1256–1264.

518  (28)  Reid, C. E.; Jerrett, M.; Petersen, M. L.; Pfister, G. G.; Morefield, P. E.; Tager, I. B.; Raffuse, S. M.;
519        Balmes, J. R. Spatiotemporal prediction of fine particulate matter during the 2008 northern
520        California wildfires using machine learning. *Environ. Sci. Technol.* **2015**, *49*, 3887–3896.

521  (29)  Liu, Y.; Paciorek, C. J.; Koutrakis, P. Estimating regional spatial and temporal variability of PM(2.5)
522        concentrations using satellite data, meteorology, and land use information. *Environ. Health*
523        *Perspect.* **2009**, *117*, 886–892.

524  (30)  Mhawish, A.; Banerjee, T.; Sorek-Hamer, M.; Lyapustin, A.; Broday, D. M.; Chatfield, R.
525        Comparison and evaluation of MODIS Multi-angle Implementation of Atmospheric Correction
526        (MAIAC) aerosol product over South Asia. *Remote Sens. Environ.* **2019**, *224*, 12–28.

527  (31)  Jethva, H.; Torres, O.; Yoshida, Y. Accuracy assessment of MODIS land aerosol optical thickness
528        algorithms using AERONET measurements over North America. *Atmos. Meas. Tech.* **2019**, *12*,
529        4291–4307.

530  (32)  Lyapustin, A.; Wang, Y.; Korkin, S.; Huang, D. MODIS Collection 6 MAIAC algorithm. *Atmos. Meas.*
531        *Tech.* **2018**, *11*, 5741–5765.

532  (33)  Hu, X.; Waller, L. A.; Lyapustin, A.; Wang, Y.; Al-Hamdan, M. Z.; Crosson, W. L.; Estes, M. G.;
533        Estes, S. M.; Quattrochi, D. A.; Puttaswamy, S. J. and Liu, Y. Estimating ground-level PM2.5
534        concentrations in the Southeastern United States using MAIAC AOD retrievals and a two-stage
535        model. *Remote Sens. Environ.* **2014**, *140*, 220–232.

536  (34)  Kloog, I.; Chudnovsky, A. A.; Just, A. C.; Nordio, F.; Koutrakis, P.; Coull, B. A.; Lyapustin, A.; Wang,
537        Y.; Schwartz, J. A New Hybrid Spatio-Temporal Model For Estimating Daily Multi-Year PM2.5
538        Concentrations Across Northeastern USA Using High Resolution Aerosol Optical Depth Data.
539        *Atmos. Environ.* **2014**, *95*, 581–590.

540  (35)  Kloog, I.; Sorek-Hamer, M.; Lyapustin, A.; Coull, B.; Wang, Y.; Just, A. C.; Schwartz, J.; Broday, D.
541        M. Estimating daily PM 2.5 and PM 10 across the complex geo-climate region of Israel using
542        MAIAC satellite-based AOD data. *Atmos. Environ.* **2015**, *122*, 409–416.

543  (36)  Just, A. C.; Wright, R. O.; Schwartz, J.; Coull, B. A.; Baccarelli, A. A.; Tellez-Rojo, M. M.; Moody, E.;
544        Wang, Y.; Lyapustin, A.; Kloog, I. Using High-Resolution Satellite Aerosol Optical Depth To

545      Estimate Daily PM2.5 Geographical Distribution in Mexico City. *Environ. Sci. Technol.* **2015**, *49*,
546      8576–8584.

547  (37)  Lee, H. J.; Chatfield, R. B.; Strawa, A. W. Enhancing the applicability of satellite remote sensing
548      for PM2.5 estimation using MODIS deep blue AOD and land use regression in california, united
549      states. *Environ. Sci. Technol.* **2016**, *50*, 6546–6555.

550  (38)  Stafoggia, M.; Bellander, T.; Bucci, S.; Davoli, M.; de Hoogh, K.; De Donato, F.; Gariazzo, C.;
551      Lyapustin, A.; Michelozzi, P.; Renzi, M.; Scortichini, M. Estimation of daily PM10 and PM2.5
552      concentrations in Italy, 2013-2015, using a spatiotemporal land-use random-forest model.
553      *Environ. Int.* **2019**, *124*, 170–179.

554  (39)  Shtein, A.; Karnieli, A.; Katra, I.; Raz, R.; Levy, I.; Lyapustin, A.; Dorman, M.; Broday, D. M.; Kloog,
555      I. Estimating daily and intra-daily PM10 and PM2.5 in Israel using a spatio-temporal hybrid
556      modeling approach. *Atmos. Environ.* **2018**, *191*, 142–152.

557  (40)  Shtein, A.; Kloog, I.; Schwartz, J.; Silibello, C.; Michelozzi, P.; Gariazzo, C.; Viegi, G.; Forastiere, F.;
558      Karnieli, A.; Just, A. C.; Stafoggia, M. Estimating Daily PM2.5 and PM10 over Italy Using an
559      Ensemble Model. *Environ. Sci. Technol.* **2019**, *54*, 120–128.

560  (41)  Zheng, C.; Zhao, C.; Zhu, Y.; Wang, Y.; Shi, X.; Wu, X.; Chen, T.; Wu, F.; Qiu, Y. Analysis of
561      influential factors for the relationship between $PM_{2.5}$ and AOD in Beijing. *Atmos. Chem. Phys.*
562      **2017**, *17*, 13473–13489.

563  (42)  Stirnberg, R.; Cermak, J.; Andersen, H. An Analysis of Factors Influencing the Relationship
564      between Satellite-Derived AOD and Ground-Level PM10. *Remote Sens (Basel)* **2018**, *10*, 1353.

565  (43)  Zhang, X.; Chu, Y.; Wang, Y.; Zhang, K. Predicting daily PM2.5 concentrations in Texas using high-
566      resolution satellite aerosol optical depth. *Sci. Total Environ.* **2018**, *631-632*, 904–911.

567  (44)  Xie, Y.; Wang, Y.; Bilal, M.; Dong, W. Mapping daily PM2.5 at 500 m resolution over Beijing with
568      improved hazy day performance. *Science of The Total Environment* **2019**, *659*, 410–418.

569  (45)  Chatfield, R. B.; Sorek-Hamer, M.; Esswein, R. F.; Lyapustin, A. Satellite mapping of $PM_{2.5}$
570      episodes in the wintertime San Joaquin Valley: a "static" model using column water vapor.
571      *Atmos. Chem. Phys.* **2020**, *20*, 4379–4397.

572  (46)  Song, W.; Jia, H.; Huang, J.; Zhang, Y. A satellite-based geographically weighted regression model
573      for regional PM2.5 estimation over the Pearl River Delta region in China. *Remote Sens. Environ.*
574      **2014**, *154*, 1–7.

575  (47)  Ma, Z.; Hu, X.; Huang, L.; Bi, J.; Liu, Y. Estimating ground-level PM2.5 in China using satellite
576      remote sensing. *Environ. Sci. Technol.* **2014**, *48*, 7436–7444.

577  (48)  Guo, Y.; Tang, Q.; Gong, D.-Y.; Zhang, Z. Estimating ground-level PM2.5 concentrations in Beijing
578      using a satellite-based geographically and temporally weighted regression model. *Remote Sens.*
579      *Environ.* **2017**, *198*, 140–149.

580    (49)    Sorek-Hamer, M.; Strawa, A. W.; Chatfield, R. B.; Esswein, R.; Cohen, A.; Broday, D. M. Improved
581            retrieval of PM2.5 from satellite data products using non-linear methods. *Environ. Pollut.* **2013**,
582            *182*, 417–423.

583    (50)    You, W.; Zang, Z.; Zhang, L.; Zhang, M.; Pan, X.; Li, Y. A nonlinear model for estimating ground-
584            level PM10 concentration in Xi'an using MODIS aerosol optical depth retrieval. *Atmos. Res.* **2016**,
585            *168*, 169–179.

586    (51)    Li, L.; Chen, B.; Zhang, Y.; Zhao, Y.; Xian, Y.; Xu, G.; Zhang, H.; Guo, L. Retrieval of daily PM2.5
587            concentrations using nonlinear methods: A case study of the beijing–tianjin–hebei region, china.
588            *Remote Sens (Basel)* **2018**, *10*, 2006.

589    (52)    He, Q.; Huang, B. Satellite-based mapping of daily high-resolution ground PM 2.5' ' in China via
590            space-time regression modeling. *Remote Sens. Environ.* **2018**, *206*, 72–83.

591    (53)    Wei, J.; Huang, W.; Li, Z.; Xue, W.; Peng, Y.; Sun, L.; Cribb, M. Estimating 1-km-resolution PM2.5
592            concentrations across China using the space-time random forest approach. *Remote Sens.*
593            *Environ.* **2019**, *231*, 111221.

594    (54)    Chen, M.-J.; Yang, P.-H.; Hsieh, M.-T.; Yeh, C.-H.; Huang, C.-H.; Yang, C.-M.; Lin, G.-M. Machine
595            learning to relate PM2.5 and PM10 concentrations to outpatient visits for upper respiratory tract
596            infections in Taiwan: A nationwide analysis. *World J Clin Cases* **2018**, *6*, 200–206.

597    (55)    Hu, X.; Belle, J. H.; Meng, X.; Wildani, A.; Waller, L. A.; Strickland, M. J.; Liu, Y. Estimating $PM_{2.5}$
598            Concentrations in the Conterminous United States Using the Random Forest Approach. *Environ.*
599            *Sci. Technol.* **2017**, *51*, 6936–6944.

600    (56)    Lary, D. J.; Lary, T.; Sattler, B. Using machine learning to estimate global PM2.5 for
601            environmental health studies. *Environ. Health Insights* **2015**, *9*, 41–52.

602    (57)    Matsuki, K.; Kuperman, V.; Van Dyke, J. A. The Random Forests statistical technique: An
603            examination of its value for the study of reading. *Sci. Stud. Read.* **2016**, *20*, 20–33.

604    (58)    Chowdhury, S.; Dey, S.; Di Girolamo, L.; Smith, K. R.; Pillarisetti, A.; Lyapustin, A. Tracking
605            ambient PM2.5 build-up in Delhi national capital region during the dry season over 15 years
606            using a high-resolution (1 km) satellite aerosol dataset. *Atmos. Environ.* **2019**, *204*, 142–150.

607    (59)    Mandal, S.; Madhipatla, K. K.; Guttikunda, S.; Kloog, I.; Prabhakaran, D.; Schwartz, J. D. Ensemble
608            averaging based assessment of spatiotemporal variations in ambient PM2.5 concentrations over
609            Delhi, India, during 2010–2016. *Atmos. Environ.* **2020**, *224*, 117309.

610    (60)    Mhawish, A.; Banerjee, T.; Broday, D. M.; Misra, A.; Tripathi, S. N. Evaluation of MODIS Collection
611            6 aerosol retrieval algorithms over Indo-Gangetic Plain: Implications of aerosols types and mass
612            loading. *Remote Sens. Environ.* **2017**, *201*, 297–313.

613    (61)    Kumar, M.; Parmar, K. S.; Kumar, D. B.; Mhawish, A.; Broday, D. M.; Mall, R. K.; Banerjee, T. Long-
614            term aerosol climatology over Indo-Gangetic Plain: Trend, prediction and potential source fields.
615            *Atmos. Environ.* **2018**, *180*, 37–50.

616 (62) Vinjamuri, K. S.; Mhawish, A.; Banerjee, T.; Sorek-Hamer, M.; Broday, D. M.; Mall, R. K.; Latif, M.
617      T. Vertical distribution of smoke aerosols over upper Indo-Gangetic Plain. *Environ. Pollut.* **2020**,
618      *257*, 113377.

619 (63) Levy, R. C.; Mattoo, S.; Munchak, L. A.; Remer, L. A.; Sayer, A. M.; Patadia, F.; Hsu, N. C. The
620      Collection 6 MODIS aerosol products over land and ocean. *Atmos. Meas. Tech.* **2013**, *6*, 2989–
621      3034.

622 (64) Hsu, N. C.; Jeong, M. J.; Bettenhausen, C.; Sayer, A. M.; Hansell, R.; Seftor, C. S.; Huang, J.; Tsay,
623      S. C. Enhanced Deep Blue aerosol retrieval algorithm: The second generation. *J. Geophys. Res.*
624      *Atmos.* **2013**, *118*, 9296–9315.

625 (65) Yuval; Levi, Y.; Dayan, U.; Levy, I.; Broday, D. M. On the association between characteristics of
626      the atmospheric boundary layer and air pollution concentrations. *Atmos. Res.* **2020**, *231*,
627      104675.

628 (66) Shafran-Nathan, R.; Yuval; Broday, D. M. Impacts of personal mobility and diurnal concentration
629      variability on exposure misclassification to ambient pollutants. *Environ. Sci. Technol.* **2018**, *52*,
630      3520–3526.

631 (67) The international centre for tropical agriculture (CIAT). *PANS Pest Articles & News Summaries*
632      **1971**, *17*, 277–279.

633 (68) Lee, H. J.; Liu, Y.; Coull, B. A.; Schwartz, J.; Koutrakis, P. A novel calibration approach of MODIS
634      AOD data to predict PM$_{2.5}$ concentrations. *Atmos. Chem. Phys.* **2011**, *11*, 7991–8002.

635 (69) Liaw, A.; Wiener, M. Classification and Regression by randomForest. *R News* **2002**, *2/3*, 18–22.

636 (70) Breiman, L. Random Forests. *Springer Science and Business Media LLC* **2001**.

637 (71) Nicodemus, K. K.; Malley, J. D.; Strobl, C.; Ziegler, A. The behaviour of random forest
638      permutation-based variable importance measures under predictor correlation. *BMC*
639      *Bioinformatics* **2010**, *11*, 110.

640 (72) Rodríguez, J. D.; Pérez, A.; Lozano, J. A. Sensitivity analysis of kappa-fold cross validation in
641      prediction error estimation. *IEEE Trans. Pattern Anal. Mach. Intell.* **2010**, *32*, 569–575.

642 (73) Wang, J.; Martin, S. T. Satellite characterization of urban aerosols: Importance of including
643      hygroscopicity and mixing state in the retrieval algorithms. *J. Geophys. Res.* **2007**, D17203, *112*.

644 (74) Xie, Y.; Wang, Y.; Zhang, K.; Dong, W.; Lv, B.; Bai, Y. Daily Estimation of Ground-Level PM2.5
645      Concentrations over Beijing Using 3 km Resolution MODIS AOD. *Environ. Sci. Technol.* **2015**, *49*,
646      12280–12288.

647 (75) Jethva, H.; Torres, O.; Field, R. D.; Lyapustin, A.; Gautam, R.; Kayetha, V. Connecting Crop
648      Productivity, Residue Fires, and Air Quality over Northern India. *Sci. Rep.* **2019**, *9*, 16594.

649 (76) Singh, N., Banerjee, T., Raju, M.P., Deboudt, K., Sorek-Hamer, M., Singh, R.S. and Mall, R.K., 2018.
650      Aerosol chemistry, transport, and climatic implications during extreme biomass burning
651      emissions over the Indo-Gangetic Plain. Atmospheric Chemistry and Physics, 18(19), pp.14197-
652      14215.

653 **TOC graphic**

654



655

656