

1 **Evaluation of machine learning techniques with multiple remote sensing datasets**
2 **in estimating monthly concentrations of ground-level PM_{2.5}**

3

4 Authors: Yongming Xu¹, Hung Chak Ho², Man Sing Wong^{2,3}, Chengbin Deng⁴, Yuan Shi⁵, Ta-
5 Chien Chan⁶, Anders Knudby⁷

6 1. School of Remote Sensing and Geomatics Engineering, Nanjing University of Information
7 Science & Technology, Nanjing, China

8 2. Department of Land Surveying and Geo-Informatics, The Hong Kong Polytechnic
9 University, Kowloon, Hong Kong

10 3. Research Institute for Sustainable Urban Development, The Hong Kong Polytechnic
11 University, Hong Kong

12 4. Department of Geography, State University of New York at Binghamton, Binghamton,
13 New York, United States

14 5. School of Architecture, Chinese University of Hong Kong, New Territories, Hong Kong

15 6. Research Center for Humanities and Social Sciences, Academia Sinica, Taiwan

16 7. Department of Geography, Environment and Geomatics, University of Ottawa, Ottawa,
17 ON, Canada

18

19 Corresponding Author: Hung Chak Ho, Department of Land Surveying and Geo-Informatics,
20 Hong Kong Polytechnic University, Hong Kong

21

22

Research Highlights

- Estimation of long-term spatially-continuous monthly PM_{2.5} dataset
- Cubist outperforms other machine learning algorithms
- Several new predictors were employed to improve the estimation of PM_{2.5}
- PM_{2.5} was estimated with a CV-RMSE of 2.64 µg/m³

23 **Abstract**

24 Fine particulate matter (PM_{2.5}) has been recognized as a key air pollutant that can
25 influence population health risk, especially during extreme cases such as wildfires. Previous
26 studies have applied geospatial techniques such as land use regression to map the ground-
27 level PM_{2.5}, while some recent studies have found that Aerosol Optical Depth (AOD) derived
28 from satellite images and machine learning techniques may be two elements that can
29 improve spatiotemporal prediction. However, there has been a lack of studies evaluating use
30 of different machine learning techniques with AOD datasets for mapping PM_{2.5}, especially in
31 areas with high spatiotemporal variability of PM_{2.5}.

32 In this study, we compared the performance of eight predictive algorithms with the use
33 of multiple remote sensing datasets, including satellite-derived AOD data, for the prediction
34 of ground-level PM_{2.5} concentration. Based on the results, Cubist, random forest and
35 eXtreme Gradient Boosting were the algorithms with better performance, while Cubist was
36 the best (CV-RMSE=2.64 µg/m³, CV-R²=0.48). Variable importance analysis indicated that the
37 predictors with the highest contributions in modelling were monthly AOD and elevation.

38 In conclusion, appropriate selection of machine learning algorithms can improve ground-
39 level PM_{2.5} estimation, especially for areas with nonlinear relationships between PM_{2.5} and
40 predictors caused by complex terrain. Satellite-derived data such as AOD and land surface
41 temperature (LST) can also be substitutes for traditional datasets retrieved from weather
42 stations, especially for areas with sparse and uneven distribution of stations.

43

44 1. Introduction

45 Fine particulate matter (PM_{2.5}) is one of the major dust-related air pollutants that can
46 increase morbidity and mortality risks, especially for cardiovascular and respiratory issues
47 (Atkinson et al., 2014). In order to reduce community health risks caused by environmental
48 exposure, previous studies have commonly applied air quality data from single or a small
49 number of monitoring stations to evaluate the temporal influences of PM_{2.5} (Liu et al., 2018;
50 Ostro et al., 2014; Wang et al., 2017), and have found positive association between PM_{2.5}
51 and chronic diseases. These results have helped pinpoint air pollution as a severe
52 community health problem (Kan et al., 2012). However, sparse distribution of air quality
53 monitoring stations across large areas reduces the ability to demonstrate the actual impact
54 of PM_{2.5} on all vulnerable populations.

55 Satellite remote sensing data can provide spatially continuous estimates of aerosol
56 optical depth (AOD), providing an alternative method to map ground-level PM_{2.5} across a
57 large region. Since AOD from satellite images has complete spatial coverage and moderate
58 spatial resolution, AOD measurement can fill in data for areas that lack monitoring stations.
59 Multiple studies have been carried out to estimate PM_{2.5} from satellite-derived AOD and
60 other environmental variables (Lai et al., 2014; Saunders et al., 2014; Wu et al., 2015). Due
61 to the spatio-temporal heterogeneity of AOD-PM_{2.5} relationships, using AOD to directly
62 represent ground-level PM_{2.5} may be inappropriate, as has been reported by previous
63 studies (Lee et al., 2011; Paciorek et al., 2008). Additional environmental predictors, such as
64 geographical and meteorological variables, have also been incorporated in models to
65 improve estimation performance (Hu et al., 2013; Kloog et al., 2011; Liu et al., 2009). To

66 derive $PM_{2.5}$ from satellite-derived AOD and other predictors, various models have been
67 developed. The most commonly used models include multiple linear regression (Lai et al.,
68 2014; Liu et al., 2004; Saunders et al., 2014; Schaap et al., 2009; Yao et al., 2018a), mixed
69 effect models (Just et al., 2015; Lee et al., 2011; Zheng et al. 2016; Xie et al., 2015), chemical
70 transport models (Crouse et al., 2016; Wang & Chen, 2016; van Donkelaar et al., 2006) and
71 geographically weighted regression (Chu et al., 2015; Chu et al., 2016; He and Huang, 2018;
72 Jiang et al., 2017; Ma et al., 2014; Shi et al., 2018; Song et al., 2014; Wu et al., 2016; You et
73 al., 2016). Recently, machine learning technology, which can fit complicated non-linear
74 relationships in many dimensions, has also been employed to derive air-pollutant
75 concentrations from remote sensing data (Chen et al., 2018; Deters et al., 2017, He & Huang,
76 2018, Yao et al., 2018b). Several machine learning methods, such as artificial neural
77 networks, generalized boosting models, support vector machine and random forest, have
78 also been used to generate models for estimating $PM_{2.5}$ (Di et al., 2016; Hu et al., 2017; Reid
79 et al., 2015; Zhan et al., 2017). However, to date, studies with machine learning for
80 estimating $PM_{2.5}$ are still rare in this field.

81 In order to better understand the potential of machine learning for $PM_{2.5}$ mapping, we
82 developed an innovative approach to estimate spatial variability of $PM_{2.5}$ by using machine
83 learning techniques with multiple predictors based on Moderate Resolution Imaging
84 Spectroradiometer (MODIS) and re-analysis data. By using machine learning techniques, it
85 can better characterize non-linear relationships for estimating air pollution based on all
86 geophysical components. To enhance the ability to develop a spatiotemporal model for
87 $PM_{2.5}$ prediction, the specific objectives of this study included 1) to develop a model for

88 predicting $PM_{2.5}$ based on remote sensing data, re-analysis data and station observed air
89 quality data; 2) to evaluate the prediction performance of different statistical methods, for
90 determining the best model setting for estimating $PM_{2.5}$; and 3) to map the spatio-temporal
91 distribution of $PM_{2.5}$ based on the best model. British Columbia of Canada was selected as
92 the case of this study, because of its complex terrain and wildfire history that can
93 significantly influence air quality across the province, including $PM_{2.5}$.

94 **2. Study Area**

95 British Columbia (BC) is the westernmost province of Canada (Fig. 1), and it is
96 characterized by mountainous terrain and heavy forest cover. BC has traditionally been
97 known for its clean environment. However, due to climate change, increasing frequency of
98 wildfires has been observed in recent decades (Wildfire Management Branch, 2014; Wotton,
99 2010). Wildfires produce excessive smoke that can influence regional air quality and severely
100 affect human health (Henderson et al., 2011; McLean et al., 2015; Krstic & Henderson et al.,
101 2015). In order to minimize air pollution risk, a National Air Pollution Surveillance (NAPS)
102 system with ground-based stations has been established across the province, monitoring
103 temporal changes in air pollutants including the daily change in $PM_{2.5}$. However, due to the
104 province's sprawling territory with complex terrain and a limited number of surveillance
105 stations, station-based observation may not be able to adequately measure the $PM_{2.5}$
106 influencing all populated regions (McLean et al. 2015). The stations with data between 2001
107 and 2014 were sparsely distributed and clustered in the southern and central parts of BC.
108 Therefore, combining satellite images to monitor the spatiotemporal changes in $PM_{2.5}$ across
109 the province is essential.

110 3. Data and Methods

111 3.1 Selection of predictors for PM_{2.5} mapping

112 According to previous studies, AOD has strong positive relationships with ground-level
113 PM_{2.5} concentrations (Engel-Cox et al., 2004; Mukai et al., 2006; Wang & Christopher, 2003;
114 Xin et al., 2014), and some studies have applied satellite-derived AOD to map PM_{2.5} (Chu et
115 al., 2016). Therefore, AOD was the first predictor for PM_{2.5} mapping. In this study, AOD data
116 were retrieved from MOD04_3K, a 3-km near-real-time aerosol dataset derived from
117 TEAAR/MODIS.

118 The PM_{2.5}-AOD relationship can be a multivariate function of a wide range of
119 influencing factors (Lary et al., 2015; Natunen et al., 2010; Song et al., 2014; van Donkelaar
120 et al., 2006). For example, meteorological and geographical predictors can be the
121 parameters of co-predicting PM_{2.5} concentrations (Jiang et al., 2017; Liu et al., 2009; Ma et
122 al., 2014; Reid et al., 2015; You et al., 2016). Built on the literature, the following parameters
123 may contribute to PM_{2.5} prediction: humidity, temperature, albedo, normalized difference
124 vegetation index (NDVI), height of the planetary boundary layer (HPBL), wind speed,
125 distance to the ocean, elevation, and calendar month. Therefore, we constructed the input
126 datasets for modelling as follows.

127 Considering the bias which sparse distribution of weather stations may produce in data
128 representing spatial variations in temperature and humidity, 26855 images of MODIS land
129 surface temperature product (MOD11A1) and 44336 images of MODIS water vapor product
130 (MOD05_L2) were used as alternatives to air temperature and relative humidity for better
131 spatial representativeness. In brief, MOD11A1 is a 1-km daily land surface temperature (LST)

132 product derived from TERRA/MODIS, and MOD05_L2 is a 1-km near-real-time water vapor
133 product derived from TERRA/MODIS.

134 In addition, NDVI and albedo were derived based on MODIS products: the MODIS
135 vegetation index product (MOD13A3), a 1-km monthly vegetation index product derived
136 from TERRA/MODIS; and the MODIS albedo product (MCD43B3), a 1-km 8-day albedo
137 product derived from TERRA/MODIS and AQUA/MODIS. For the mapping purpose, all MODIS
138 datasets were re-projected to the Albers projection, resampled to 1-km spatial resolution,
139 and averaged for each month.

140 Finally, HPBL and wind speed were derived from NCAR/NCEP re-analysis data, which
141 provides the corresponding data on a monthly basis. Elevation was derived from a digital
142 elevation model (DEM) dataset of the Shuttle Radar Topography Mission (SRTM). Distance to
143 the ocean was calculated by buffer analysis based on the coastal boundary of BC.

144 Based on the satellite-derived products and re-analysis data, a total of 10 predictors
145 were employed to estimate ground-level $PM_{2.5}$ concentration across BC: monthly AOD,
146 monthly vapor, monthly LST, monthly NDVI, monthly albedo, monthly HPBL, monthly wind
147 speed, elevation, distance to ocean and calendar month (Table 1).

148 It is known that the relationship between environmental predictors and $PM_{2.5}$ may vary
149 across space (Hu et al., 2013; Song et al., 2014), as well as time. We did not include spatial
150 predictors (e.g. latitude, longitude) other than “distance to ocean”, and we did not use
151 spatially weighted models such as geographically weighted regression, because of the
152 limited insight that can be gained from using such predictors/models, and the limited

153 transferability such models will have to other geographical regions.

154 **3.2 Model development with machine learning algorithms**

155 Association between $PM_{2.5}$ concentration of air quality monitoring stations and the
156 values of predictors retrieved by the locations of stations were first established for each
157 machine learning model in order to estimate the spatial distributions of ground-level $PM_{2.5}$
158 concentrations. In this study, ground-level $PM_{2.5}$ concentrations for modelling were retrieved
159 from 63 stations of the NAPS network operated by Environment Canada, with hourly $PM_{2.5}$
160 data between 2001 and 2014 across BC. Since several stations within this study period did
161 not provide temporal-continuous observations, or even had significant data gaps in temporal
162 observation, we averaged hourly $PM_{2.5}$ data on a daily basis, then converted the daily
163 information to the monthly average $PM_{2.5}$ concentrations based on all valid daily values.

164 These monthly average $PM_{2.5}$ values across BC province were then applied to the
165 following statistic algorithms to construct the regression models: 1) multiple linear
166 regression (MLR), 2) Bayesian Regularized Neural Networks (BRNN), 3) Support Vector
167 Machines with Radial Basis Function Kernel (SVM), 4) Least Absolute Shrinkage and Selection
168 Operator (LASSO), 5) Multivariate Adaptive Regression Splines (MARS), 6) Random forest
169 (RF), 7) eXtreme Gradient Boosting (XGBoost), and 8) Cubist.

170 MLR is a widely used algorithm in remote sensing applications because of its simplicity,
171 but it relies on several assumptions concerning data distributions, and its performance
172 depends on meeting these assumptions as well as the linearity of the modeled relationship
173 (Helsel and Hirsch 1992). BRNN is a back-propagation network that based on a mathematical
174 technique named Bayesian regularization to convert nonlinear regression into “well-posed”

175 problems (Burden and Winkler, 2008). It is more robust than standard back-propagation
176 neural networks. SVM was originally developed for classification by constructing separating
177 hyperplanes to define decision boundaries, and later expanded for regression. To map
178 samples to high dimension space, kernel functions were introduced. The radial basis function
179 showed its advances of handling nonlinear problems and fewer tunable parameters (Hsu, 2003;
180 Bennett and Campbell, 2000). LASSO is a regularization and variable selection method which
181 shrinks coefficients by forcing some less important coefficients to zero (Tibshirani, 1996). It can
182 improve the model interpretability and reduce overfitting. MARS is a fully automated method
183 based on the divide-and-conquer strategy, in which the training dataset is split into
184 piecewise linear segments (splines) (Friedman, 1991). RF is an ensemble-based decision tree
185 approach, which consists of a combination of decision trees fitted by randomly selected
186 subsets of training samples. Final predictions produced by RF model are determined by the
187 average of the results of all the trees (Breiman, 2001). XGBoost is an ensemble tree method
188 which follows the principle of Gradient boosting framework (Friedman, 2001), and uses
189 regularization techniques to control overfitting and model complexity (Chen and Guestrin,
190 2016). Cubist is a rule-based tree model, which produces multiple linear regression models
191 in the terminal nodes of trees based on the M5 theory (Quinlan, 1992; RuleQuest, 2018). A
192 prediction at the terminal node is made by the corresponding linear regression model and is
193 smoothed by combining with predictions from nearest-neighbor nodes within the tree to
194 improve prediction accuracy (Houborg & McCabe, 2018). In addition, Cubist also constructs
195 multiple tree models (called committees), each of which consists of a set of rule-based
196 models (John et al., 2018). Predictions from all the committees are averaged to produce the

197 final prediction.

198 Except for the widely-used traditional MLR algorithm, others were machine learning
199 algorithms, which can effectively fit nonlinear and complex relationships between outcomes
200 and predictors (Ngufor et al. 2015). In this study, the complex terrain of the study area can
201 form a nonlinear relationship between ground-level PM_{2.5} concentrations and all predictors,
202 for which machine learning models may provide better results.

203 In order to optimize the PM_{2.5} estimation, parameter values were adjusted in each
204 machine learning model with a fitting process, based on the determination of the best
205 parameters by cyclic testing with committees of 1, 5, 10, 20, 50, and neighbors of 0, 1, 5, 9.
206 In addition, predictions of PM_{2.5} concentrations with all machine learning models were
207 conducted with the R (R Development Core Team).

208 **3.3 Model evaluation**

209 10-fold cross-validation was performed to evaluate the accuracy of all machine learning
210 models. Data were first randomly divided into 10 subsets, with one of the subsets used as
211 the validation dataset and the remaining used as training datasets; then repeating 10 times
212 until all subsets have been used as validation datasets once. Root-mean-square error (CV-
213 RMSE) and coefficient of determination (CV-R²) based on the comparison of validation and
214 training data were used to evaluate the accuracy of each machine learning model. While the
215 best model for PM_{2.5} estimation was determined based on the accuracies, variable
216 importance analysis was also conducted to evaluate the contributions of each predictor in
217 PM_{2.5} estimation, based on the determination of percentage increase in mean square error
218 (%IncMSE) of each model relative to the original error, after a predictor was randomly

219 permuted. A higher value of %IncMSE indicated higher importance of this corresponding
220 predictor to the estimation.

221 **4. Results**

222 **4.1 Empirical relationship between PM_{2.5} and AOD**

223 A total of 1242 records of observed data of ground-level PM_{2.5} concentrations were
224 retrieved from stations with effective monthly AOD values based on location. In brief, PM_{2.5}
225 concentrations of this subset ranged from 1.26µg/m³ to 51.14µg/m³, with an average of
226 5.26µg/m³ and a median of 4.58µg/m³. This indicated a clean environment with low air
227 pollution during the study period across BC, except in a few extreme cases. Based on the
228 observed data, the extremes in PM_{2.5} concentration samples were observed in August 2003
229 and August 2010, when there were wildfire events (e.g. 2003 Okanagan Mountain Park Fire)
230 across BC.

231 A positive but poor correlation was observed based on evaluation of an empirical
232 relationship between observed PM_{2.5} and satellite-derived AOD (Fig. 2), with a correlation
233 coefficient (R) of 0.34 (P-value < 0.01), a clustering of data was found with AOD value less
234 than 0.8 and PM_{2.5} value less than 15µg/m³. Observed data with moderate or high values
235 were scattered, possibly due to the complexity of the atmospheric conditions and
236 landscapes across BC. Similar evidence has also been found in a previous study, which
237 demonstrated a non-linear relationship between geophysical environment and air
238 temperature across BC (Xu et al., 2014). Therefore, the use of simple linear regression for
239 ground-level PM_{2.5} estimation is insufficient and inaccurate, and nonlinear multivariate
240 models should be adopted to predict PM_{2.5} under consideration of relevant atmosphere-

241 surface interactions.

242 **4.2 Model performance**

243 Parameters of machine learning models were optimized with the fitting process, by
244 cyclic testing with a given parameter range and step size. Based on the results of optimized
245 models, CV-RMSE ranged from 2.64 $\mu\text{g}/\text{m}^3$ to 3.24 $\mu\text{g}/\text{m}^3$ and CV-R² ranged from 0.22-0.49
246 (Table 2). Among all, RF, XGBoost and Cubist were the models with better performance,
247 while Cubist had the best performance determined by CV-RMSE. With 20 committees and 5
248 neighbors as optimal parameters, CV-RMSE and CV-R² of Cubist were 2.64 $\mu\text{g}/\text{m}^3$ and 0.48. In
249 contrast, MLR method had the lowest performance (CV-RMSE=3.24 $\mu\text{g}/\text{m}^3$ and CV-R²=0.22),
250 indicating its poor capability of capturing complex relationships for the study area.

251 For the best model, the predicted and observed values were well aligned with the line
252 of best fit (Fig. 3), indicating the high accuracy of PM_{2.5} estimation with Cubist. However,
253 underestimation was also found for observed data with high PM_{2.5} values (> 20 $\mu\text{g}/\text{m}^3$),
254 possibly due to the small sample size, resulting in inability to robustly predict these high-
255 value data with a decision-based machine learning algorithm. Moreover, average deviation
256 of PM_{2.5} estimation was 0.07 $\mu\text{g}/\text{m}^3$, slightly higher than the deviation of observed values.
257 These results show that lower PM_{2.5} concentration in observed data may result in
258 overestimation, while higher values in observed data might result in underestimation during
259 prediction.

260 **4.3 Variable importance analysis**

261 Based on the variable importance analysis, the predictors with highest contributions to
262 the Cubist model were monthly AOD and elevation. %IncRMSE without monthly AOD as

263 predictor was 12.14%, possibly due to its strong association between AOD and ground-level
264 air quality. %IncRMSE without elevation as a predictor was 9.26%, also suggesting a high
265 importance in PM_{2.5} estimation because of the influences of complex terrain in BC, with
266 great variations in altitude between the coast and interior. However, there shall be several
267 factors which contributed to the importance of elevation for predictions of PM_{2.5}: areas with
268 high elevation are inclined to suffer from wildfires; areas with low elevation tend to be
269 influenced by human activities. As AOD is an important predictor in the models, elevation
270 may be used to correct for model predictions. In addition, %IncRMSE of monthly albedo,
271 monthly LST and calendar month ranged from 4% to 6%. Predictors with the least
272 importance were monthly wind speed, monthly HPBL, monthly vapor and monthly NDVI,
273 with a range of %IncRMSE between 2% and 4%.

274 **4.4 Determination of location-based error**

275 To further determine the spatial variability of error, RMSEs were extracted by the
276 location of each station (Fig. 4). Most stations had RMSEs lower than 2.0µg/m³, while the
277 stations with the lowest RMSEs were in southeastern, western and southwestern BC. In
278 contrast, high errors were found at stations located in central and central-southern parts of
279 BC, with RMSEs ranging from 3.0 - 4.0µg/m³ or even higher. Compared with the DEM, these
280 stations with higher RMSEs were in mountainous valleys with high PM_{2.5} concentrations.
281 Estimation errors of these stations were mostly negative, indicating an underestimation of
282 ground-level PM_{2.5} across these valleys. These were also aligned with previous findings (Fig.
283 3) that observed data with higher PM_{2.5} may introduce a higher chance of underestimation
284 based on the Cubist model in this study.

285 5 Discussion

286 5.1 Spatiotemporal variability of ground-level PM_{2.5} concentration

287 Based on the average concentrations of ground-level PM_{2.5} between 2001 and 2014
288 (Fig. 6), considerable spatial heterogeneity was found across BC. Generally, northern and
289 northeastern BC were areas with lower PM_{2.5} concentrations (< 4 µg/m³), while mountainous
290 regions across western BC were areas with higher concentrations of PM_{2.5} (5-6 µg/m³). We
291 also observed several extreme cases in mountainous valleys of BC (>7 µg/m³). One reason
292 for this spatiotemporal variability might be associated with wildfires, as this was a major
293 source of ambient PM_{2.5} across mountainous BC. Previous studies have found a particular
294 deposition process of PM_{2.5}, emitted from biomass burning, with long-distance transport
295 (Ward et al., 1991; Sapkota et al., 2005). We should emphasize that terrain can play an
296 influential role in the deposition, due to the aerodynamic characteristics of PM_{2.5} and the
297 topographical effect on wind flow. For example, the mountainous topography of BC, with its
298 irregular terrain, can result in uneven distribution of air pressure that further influences
299 near-surface wind. The effect of local terrain on PM_{2.5} dispersion due to its impact on wind
300 dynamics has also been found in another study in mountainous areas (Shi et al., 2017). A
301 considerable fraction of PM_{2.5} is therefore expected to be trapped by the leeward side of
302 mountains, valleys, canyons and basins (Steyn et al., 2013) under the typical transport
303 process of air pollutants. Urban areas with high aerodynamic surface roughness may also
304 have influence similar to this topographical effect on the deposition of PM_{2.5} from wildfires
305 (Landsberg, 1981). These findings indicate that regions across BC with lower altitude and
306 with poorer air dispersion due to topographical effects may be areas with higher PM_{2.5}

307 concentration. In addition, these facts may also partly explain the lower contribution of
308 monthly coarse spatial resolution (2.5 degree latitude x 2.5 degree longitude) and monthly
309 wind speed in modelling based on variable importance analysis, while another reason may
310 be the coarse spatial resolution (2.5 degree) of predictors derived from NCEP/NCAR re-
311 analysis data. Due to this resolution, it cannot represent micro-scale topographical effects
312 on air pollution transport and deposition. Some mountain valleys in BC have high
313 temperatures and little rainfall during the summer, and become dry enough to have near-
314 desert conditions with substantial amounts of dust suspended in the atmosphere, which is
315 also contributed to the high PM_{2.5} concentrations of valleys. An isolated cluster of high
316 PM_{2.5} in the Greater Vancouver Area and its surrounding regions was also observed, which
317 has not been shown in other BC cities. This can be attributed to the large population and
318 corresponding industrial, traffic and domestic emissions over this region.

319 Furthermore, CV-RMSE of this study was lower than previous research in other areas
320 (Liu et al., 2009; Song et al., 2014; Kloog et al., 2014; You et al., 2015; Reid et al., 2015; Liu et
321 al., 2005), partially indicating better air quality of BC compared to other regions. In contrast,
322 a lower CV-R² was found, which may be the result of extreme wildfire events in BC leading to
323 data with high PM_{2.5} concentration values as outliers in modelling.

324 **5.2 Advantages and Limitations**

325 In this study, optimization of machine learning models can effectively reduce the
326 sensitivity of the model tree to data noise with uncertainty; while the evaluation of eight
327 machine learning algorithms for modelling indicated that ensemble machine learning can
328 improve the accuracy of ground-level PM_{2.5} prediction. In addition, weather stations were

329 generally designed under government protocols, resulting in a sparse and uneven
330 distribution. This, as well as the strong variation in topography across the study area, makes
331 it unsuitable to apply conventional geostatistical methods such as spatial interpolation for
332 mapping the spatial variability of environmental variables (e.g. temperature and humidity),
333 while these maps should be the input layers for air quality prediction. In this study, we
334 provided an alternative, in which the use of LST and atmospheric water vapor derived from
335 satellite images can be substitutes for temperature and humidity maps.

336 There were areas with data missing from the prediction (Fig. 6). These were mainly the
337 high-altitude areas covered with perennial snow, because the Dark Target algorithm for AOD
338 retrieval was designed for areas with lower surface reflectance under a clear sky. For areas
339 with high surface reflectance values (e.g. snow coverage and desert), null values of AOD data
340 would be found. In addition, AOD values surrounding the missing data were generally high,
341 because AOD in such areas could be easily overestimated by the Dark Target algorithm,
342 especially in areas with high surface brightness and low vegetation coverage (Levy et al.,
343 2010). These became the areas with missing values of $PM_{2.5}$ concentration across snow
344 coverage in this study, and there were extremely high values of $PM_{2.5}$ concentration
345 surrounding these areas with missing data, especially those areas just below the snowline
346 with lower vegetation coverage. The issue of missing data is especially noticeable in winter,
347 as mountainous BC was covered by snow, resulting in high surface reflectance, and this area
348 was also constantly covered by clouds due to the relatively humid weather in wintertime,
349 resulting in spatiotemporal incompleteness of $PM_{2.5}$ estimation.

350 In addition, the $PM_{2.5}$ concentration over BC showed high values both in western high

351 mountains and the Fraser River Delta. The principal sources of $PM_{2.5}$ is likely different
352 between these areas. In mountain areas high $PM_{2.5}$ concentration is mostly caused by
353 wildfires, while in the Fraser River Delta high $PM_{2.5}$ concentration is caused by human
354 activity. Due to the lack of the chemical characteristics of particulate matter, we cannot
355 perform a chemical analysis of fine particulate matter over these regions. Further study with
356 field measurement should be applied to observe personal and ambient exposure of $PM_{2.5}$
357 from multiple sources. However, this future study will be limited by the accessibility of field
358 measurement and the potential bias from indoor-outdoor exchange of air pollution.

359 **6 Conclusions**

360 In this study, we evaluated the abilities of machine learning techniques to estimate the
361 monthly concentrations of ground-level $PM_{2.5}$ between 2001 and 2014, based on eight
362 algorithms with predictors derived from remote sensing and meteorological re-analysis data.
363 Predictions from these algorithms were evaluated by a 10-fold cross-validation, with CV-
364 RMSE ranging from $2.64\mu\text{g}/\text{m}^3$ to $3.25\mu\text{g}/\text{m}^3$ and CV- R^2 ranging from 0.23-0.49. Among all,
365 Cubist had the best performance (CV-RMSE= $2.64\mu\text{g}/\text{m}^3$, CV- $R^2=0.48$). A series of maps were
366 produced for representing the monthly $PM_{2.5}$ concentrations across BC, which can be
367 reference information on intra-province air pollution over 14 years for further air quality
368 monitoring and public health surveillance. In conclusion, selection of appropriate machine
369 learning algorithms for modelling can improve the accuracy in $PM_{2.5}$ estimation, while using
370 satellite-derived data as predictors can minimize the spatial bias compared with use of
371 traditional datasets retrieved from weather stations.

372 Recently, deep learning technology has attracted much attention in various fields.

373 Compared with conventional machine learning technology, deep learning can provide better
374 accuracy but requires a large amount of training data (Camilleri and Prescott, 2017; Ravi et
375 al., 2017). Due to the limited number of air quality stations, there are not enough samples to
376 sufficiently train deep learning models. Therefore it is a big challenge to adopt deep learning
377 technology to map $PM_{2.5}$ at the present stage. In the future, if the big training data
378 requirement of deep learning can be resolved, it is expected to achieve improved estimation
379 of $PM_{2.5}$ concentration from remote sensing data. The method used in this study with the
380 combination of machine learning and multi-source variables was a preliminary attempt to
381 map $PM_{2.5}$ concentration with the currently available data and suitable machine learning
382 methods. The method proposed in this paper could also be applied to other complex terrain
383 regions with sparse distributed air quality stations. Due to the limitation of AOD retrieval
384 algorithms, the remotely sensed AOD data have coarse spatial resolutions. Re-analysis data have
385 even coarser resolutions. The low spatial resolution of datasets restricts the application of this
386 method on a small scale (e.g. city scale).

387

388 **Acknowledgments**

389 This work was supported by the Social Sciences Foundation of the Ministry of Education
390 of China (Grant No. 17YJCZH205) and the National Key Research and Development Program
391 of China (2017YFB0503903-4). We would like to thank the Land Processes Distributed Active
392 Archive Center (LPDAAC) and Level-1 and Atmosphere Archive & Distribution System
393 (LAADS) for providing MODIS data, US Geological Survey (USGS) for providing SRTM/DEM
394 data, and National Oceanic and Atmospheric Administration (NOAA)/ Earth System Research

395 Laboratory (ESRL) for providing NCEP Reanalysis data. Man Sing Wong thanks the support in
396 part by a grant from the General Research Fund (project ID: 15205515); and a grant of PolyU
397 1-ZVFD from the Research Institute for Sustainable Urban Development, the Hong Kong
398 Polytechnic University. We also thank the two reviewers for their valuable comments and
399 suggestions.

400

401 **References**

402 Atkinson, R. W., Kang, S., Anderson, H.R., Mills, I.C., Walton, H.A., 2014. Epidemiological time
403 series studies of PM2.5 and daily mortality and hospital admissions: a systematic review
404 and meta-analysis. *Thorax* 69, 660–665

405 B.C. Wildfire Management Branch. 2014. Proactive Wildfire Threat Reduction. Accessed June
406 15, 2017. [http://docs.openinfo.gov.bc.ca/d63519414a_response_package_fnr-2014-](http://docs.openinfo.gov.bc.ca/d63519414a_response_package_fnr-2014-00274.pdf)
407 [00274.pdf](http://docs.openinfo.gov.bc.ca/d63519414a_response_package_fnr-2014-00274.pdf)

408 Bennett, K.P., Campbell, C., 2000. Support vector machines: hype or hallelujah? *SIGKDD*
409 *Explor.* 2, 1–13.

410 Breiman, L., 2001. Random forests. *Mach. Learn.* 45, 5–32.

411 Burden, F., Winkler, D., 2008. Bayesian regularization of neural networks. *Methods Mol. Biol.*
412 458, 25–44.

413 Camilleri, D. Prescott, T., 2017. Analysing the limitations of deep learning for developmental
414 robotics. In: *Biomimetic and Biohybrid Systems. 6th International Conference, Living*
415 *Machines 2017, Stanford, CA, USA.*

416 Chen, B., Song, Y., Jiang, T., Chen, Z., Huang, B., Xu, B., 2018. Real-time estimation of

417 population exposure to PM2.5 using mobile-and station-based big data. *Int. J. Environ.*
418 *Res. Public Health* 15, 573.

419 Chen, T., Guestrin, C., 2016. XGBoost: A scalable tree boosting system. *Proceedings of the*
420 *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining,*
421 *785-789.*

422 Chu, H.J., Huang, B., Lin, C.Y., 2015. Modeling the spatio-temporal heterogeneity in the
423 PM10-PM2.5 relationship. *Atmos. Environ.* 102, 176–182.

424 Chu, Y., Liu, Y., Li, X., Liu, Z., Lu, H., Lu, Y., Mao, Z., Chen, X., Li, N., Ren, M., Liu, F., Tian, L.,
425 Zhu, Z., Xiang, H., 2016. A review on predicting ground PM2.5 concentration using
426 satellite aerosol optical depth. *Atmosphere* 7, 129.

427 Crouse, D.L., Philip, S., van Donkelaar, A., Martin, R.V., Jessiman, B., Peters, P.A.,
428 Weichenthal, S., Brook, J.R., Hubbell, B., Burnett, R.T., 2016. A new method to jointly
429 estimate the mortality risk of long-term exposure to fine particulate matter and its
430 components. *Sci. Rep.* 6, 18916.

431 Deters, J.K., Zalakeviciute, R., Gonzalez, M., Rybarczyk, Y., 2017. Modeling PM2.5 urban
432 pollution using machine learning and selected meteorological parameters. *J. Elect.*
433 *Comput. Eng.* 2017, 1-14

434 Di, Q., Koutrakis, P., Schwartz, J., 2016. A hybrid prediction model for PM2.5 mass and
435 components using a chemical transport model and land use regression. *Atmos. Environ.*
436 131, 390–399.

437 Engel-Cox, J.A., Holloman, C.H., Coutant, B.W., Hoff, R.M., 2004. Qualitative and quantitative
438 evaluation of MODIS satellite sensor data for regional and urban scale air quality.

439 Atmos. Environ. 38, 2495–2509.

440 Friedman, J.H., 1991. Multivariate Adaptive Regression Splines. *Ann. Stat.* 19, 1–67.

441 Friedman, J.H., 2001. Greedy function approximation: a gradient boosting machine. *Ann.*
442 *Stat.* 29, 1189–1232.

443 He, Q., Huang, B., 2018. Satellite-based high-resolution PM_{2.5} estimation over the Beijing-
444 Tianjin-Hebei region of China using an improved geographically and temporally
445 weighted regression model. *Environ. Pollut.* 236, 1027–1037.

446 Helsel, D. R., Hirsch, R. M., 1992. *Statistical Methods in Water Resources*, 296–299.
447 Amsterdam: Elsevier.

448 Henderson, S.B., Brauer, M., MacNab, Y.C., Kennedy, S.M., 2011. Three measures of forest
449 fire smoke exposure and their associations with respiratory and cardiovascular health
450 outcomes in a population-based cohort. *Environ. Health Perspect.* 119, 1266–1271.

451 Houborg, R., McCabe, M.F., 2018. A hybrid training approach for leaf area index estimation
452 via Cubist and random forests machine-learning. *ISPRS J. Photogramm. Remote Sens.*
453 135, 173–188.

454 Hsu, C.W., Chang, C.C., Lin, C.J., 2003. A practical guide to support vector classification.

455 Hu, X., Waller, L.A., Al-Hamdan, M.Z., Crosson, W.L., Estes, M.G., Jr, Estes, S.M., Quattrochi,
456 D.A., Sarnat, J.A., Liu, Y., 2013. Estimating ground-level PM_{2.5} concentrations in the
457 southeastern U.S. using geographically weighted regression. *Environ. Res.* 121, 1–10.

458 Hu, X., Belle, J.H., Xia, M., Wildani, A., Waller, L., Strickland, M., Liu Y., 2017. Estimating
459 pm_{2.5} concentrations in the conterminous United States using the random forest
460 approach. *Environ. Sci. Technol.* 51, 6936–6944.

461 Jiang M., Sun W., Yang G., Zhang D., 2017. Modelling seasonal GWR of daily PM2.5 with
462 proper auxiliary variables for the Yangtze River Delta. *Remote Sens.* 9, 346.

463 John, R., Chen, J., Giannico, V., Park, H., Xiao, J., Shirkey, G., Ouyang, Z., Shao G., Laforteza,
464 R., Qi, J., 2018. Grassland canopy cover and aboveground biomass in Mongolia and
465 Inner Mongolia: Spatiotemporal estimates and controlling factors. *Remote Sens.*
466 *Environ.* 213, 34–48.

467 Kan, H., Chen, R., Tong, S., 2012. Ambient air pollution, climate change, and population
468 health in China. *Environ. Int.* 42, 10–19.

469 Kloog, I., Sorek-Hamer, M., Lyapustin, A., Coull, B., Wang, Y., Just, A. C., Schwartz, J., Broday,
470 D. M., 2015. Estimating daily pm 2.5, and pm 10, across the complex geo-climate region
471 of Israel using MAIAC satellite-based AOD data. *Atmos. Environ.* 122, 409–416.

472 Kloog, I., Koutrakis, P., Coull, B.A., Lee, H.J., Schwartz, J., 2011. Assessing temporally and
473 spatially resolved PM2.5 exposures for epidemiological studies using satellite aerosol
474 optical depth measurements. *Atmos. Environ.* 45, 6267–6275.

475 Krstic, N., Henderson, S.B., 2015. Use of MODIS data to assess atmospheric aerosol before,
476 during, and after community evacuations related to wildfire smoke. *Remote Sens.*
477 *Environ.* 166, 1–7.

478 Lai, H.K., Tsang, H., Thach, T.Q., Wong, C.M., 2014. Health impact assessment of exposure to
479 fine particulate matter based on satellite and meteorological information. *Environ. Sci.*
480 *Process. Impact* 2014, 16, 239–246.

481 Landsberg, H.E., 1981. *The urban climate* (Vol. 28). Academic Press.

482 Lary, D.J., Lay, T., Sattler, B., 2015. Using machine learning to estimate global PM2.5 for

483 environmental health studies, *Environ. Health Insights* 9, 41–52.

484 Lee, H.J., Chatfield, R.B., Strawa, A.W., 2016. Enhancing the applicability of satellite remote
485 sensing for PM_{2.5} estimation using MODIS deep blue AOD and land use regression in
486 California, United States. *Environ. Sci. Technol.* 50, 6546–6555.

487 Lee, H.J., Liu, Y., Coull, B. A., Schwartz, J., Koutrakis, P., 2011. A novel calibration approach of
488 MODIS AOD data to predict PM_{2.5} concentrations. *Atmospheric Chem. Phys.* 11, 7991–
489 8002.

490 Levy, R.C., Remer, L.A., Kleidman, R.G., Mattoo, S., 2010. Global evaluation of the collection
491 5 modis dark-target aerosol products over land. *Atmospheric Chem. Phys.*, 10, 10399–
492 10420.

493 Liu, J., Li, W., Wu, J., Liu, Y. 2018. Visualizing the intercity correlation of PM_{2.5} time series in
494 the Beijing-Tianjin-Hebei region using ground-based air quality monitoring data. *PloS*
495 *one*, 13, e0192614.

496 Liu, Y., 2014. Mapping annual mean ground-level PM_{2.5} concentrations using multiangle
497 imaging spectroradiometer aerosol optical thickness over the contiguous United States.
498 *J. Geophys. Res.* 109, D22.

499 Liu Y., Paciorek C.J., Koutrakis P., 2009. Estimating regional spatial and temporal variability of
500 PM_{2.5} concentrations using satellite data, meteorology, and land use information.
501 *Environ. Health Perspect.* 117, 886–892.

502 Liu, Y., Franklin, M., Kahn, R., Koutrakis, P., 2007. Using aerosol optical thickness to predict
503 ground-level PM 2.5 concentrations in the St. Louis area: a comparison between MISR
504 and MODIS. *Remote Sens. Environ.* 107, 33–44.

505 Ma, Z., Hu, X., Huang, L. Bi, J., Liu, Y., 2014. Estimating ground-level PM_{2.5} in China using
506 satellite remote sensing. *Environ. Sci. Technol.* 48, 7436–7444.

507 McLean, K.E., Yao, J., Henderson, S.B., 2015. An evaluation of the British Columbia Asthma
508 Monitoring System (BCAMS) and PM_{2.5} exposure metrics during the 2014 forest fire
509 season. *Int. J. Environ. Res. Public Health* 12, 6710–6724.

510 Mukai, S., Sano, I., Satoh, M., Holben, B.N., 2006. Aerosol properties and air pollutants over
511 an urban area. *Atmos. Res.* 82, 643–651.

512 Natunen, A., Arola, A., Mielonen, T., Huttunen, J., Komppula, M., Lehtinen, K.E.J., 2010. A
513 multi-year comparison of PM_{2.5} and AOD for the Helsinki region. *Boreal Environ. Res.*
514 15, 544–552

515 Ngufor, C., Murphree, D., Upadhyaya, S., Madde, N., Kor, D., Pathak, J., 2015. Effects of
516 plasma transfusion on perioperative bleeding complications: a machine learning
517 approach. *Stud. Health Technol. Inform.* 216, 721–725.

518 Ostro, B., Malig, B., Broadwin, R., Basu, R., Gold, E.B., Bromberger, J.T., Derby, C., Feinstein,
519 S., Greendale, G. Jackson, E., Kravitz, H.M., Matthews, K.A., Sternfeld, B., Tomey, K.,
520 Green, R.R., Green. R., 2014. Chronic PM_{2.5} exposure and inflammation: determining
521 sensitive subgroups in mid-life women. *Environ. Res.* 132, 168–175.

522 Paciorek, C.J., Liu, Y., Moreno-Macias, H., Kondragunta, S., 2008. Spatiotemporal
523 associations between GOES aerosol optical depth retrievals and ground-level PM_{2.5}.
524 *Environ. Sci. Technol.* 42, 5800–5806.

525 R Core Development Team, 2016. R: A language and environment for statistical computing. R
526 Foundation for Statistical Computing, Vienna, Austria.

527 Ravi, D., Wong, C., Deligianni, F., Berthelot, M., Andreu-Perez, J., Lo, B., Yang, G.Z. 2017.
528 Deep learning for health informatics. *IEEE J. Biomed. Health Inform.* 21, 4-21.

529 Reid C.E., Jerrett, M., Petersen, M.L., Pfister, G.G., Morefield, P.E., Tager, I.B., Raffuse, S.M.,
530 Balmes, J.R., 2015. Spatiotemporal prediction of fine particulate matter during the 2008
531 Northern California wildfires using machine learning. *Environ. Sci. Technol.* 49,
532 3887-3896

533 RuleQuest., 2018. Data mining with Cubist, <https://www.rulequest.com/cubist-info.html>

534 Sapkota, A., Symons, J.M., Kleissl, J., Wang, L., Parlange, M.B., Ondov, J., Breyse, P.N.,
535 Buckley, T.J., 2005. Impact of the 2002 Canadian forest fires on particulate matter air
536 quality in Baltimore City. *Environ. Sci. Technol.* 39, 24-32.

537 Saunders, R.O., Kahl, J.D.W., Ghorai, J.K., 2014. Improved estimation of PM2.5 using
538 Lagrangian satellite-measured aerosol optical depth. *Atmos. Environ.* 91, 146-153.

539 Schaap, M., Apitley, A., Timmermans, R.M.A., Koelemeijer, R.B.A., de Leeuw G., 2009.
540 Exploring the relation between aerosol optical depth and PM2.5 at Cabauw, the
541 Netherlands. *Atmos. Chem. Phys.* 9, 909-925.

542 Shi, Y., Ho, H.C., Xu, Y., Ng, E., 2018. Improving satellite aerosol optical Depth-PM2.5
543 correlations using land use regression with microscale geographic predictors in a high-
544 density urban context. *Atmos. Environ.* Doi: 10.1016/j.atmosenv.2018.07.021.

545 Shi, Y., Lau, K.K.L., Ng, E., 2017. Incorporating wind availability into land use regression
546 modelling of air quality in mountainous high-density urban environment. *Environ. Res.*
547 157, 17-29.

548 Song, W., Jia, H., Huang, J., Zhang, Y., 2014, A satellite-based geographically weighted

549 regression model for regional PM_{2.5} estimation over the Pearl River Delta region in
550 China. *Remote Sens. Environ.* 154, 1–7.

551 Steyn, D.G., De Wekker, S.F., Kossmann, M., Martilli, A., 2013. Boundary layers and air
552 quality in mountainous terrain. In *Mountain Weather Research and Forecasting*.
553 Springer Netherlands, pp. 261–289.

554 Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Series B*
555 *Stat. Methodol.* 58, 267–288.

556 van Donkelaar A., Martin R.V., Park R. J., 2006. Estimating ground-level pm_{2.5} using aerosol
557 optical depth determined from satellite remote sensing. *J. Geophys. Res.* 111, D21.

558 Wang, B., Chen, Z., 2016. High-resolution satellite-based analysis of ground-level PM_{2.5} for
559 the city of Montreal. *Sci. Total Environ.* 541, 1059–1069.

560 Wang, J., Christopher, S.A., 2003. Intercomparison between satellite derived aerosol optical
561 thickness and PM_{2.5} mass: implications for air quality studies, *Geophys. Res. Lett.* 30,
562 2095.

563 Wang, Y., Shi, L., Lee, M., Liu, P., Di, Q., Zanobetti, A., Schwartz, J.D., 2017. Long-term
564 exposure to PM_{2.5} and mortality among older adults in the Southeastern US.
565 *Epidemiology* 28, 207–214.

566 Ward, D.E., Hardy, C.C., 1991. Smoke emissions from wildland fires. *Environ. Int.* 17, 117-
567 134.

568 Xie, Y., Wang, Y., Zhang, K., Dong, W., Lv, B., Bai, Y., 2015. Daily estimation of ground-level
569 PM_{2.5} concentrations over Beijing using 3km resolution MODIS AOD. *Environ. Sci.*
570 *Technol*, 19, 12280–12288.

571 Wotton, B.M., Nock, C.A., Flannigan, M.D., 2010. Forest fire occurrence and climate change
572 in Canada. *Int. J. Wildland Fire* 19, 253–271.

573 Wu, J., Li, J., Peng, J., Li, W., Xu, G., Dong, C., 2015. Applying land use regression model to
574 estimate spatial variation of PM_{2.5} in Beijing, China. *Environ. Sci. Pollut. Res. Int.* 22,
575 7045-7061.

576 Wu, J., Yao, F., Li, W., Si, M., 2016. VIIRS-based remote sensing estimation of ground-level
577 PM_{2.5} concentrations in Beijing-Tianjin-Hebei: A spatiotemporal statistical model.
578 *Remote Sens. Environ.*, 184, 316–328.

579 Xin, J., Zhang, Q., Wang, L., Gong, C., Wang, Y., Liu, Z., Gao, W., 2014. The empirical
580 relationship between the PM_{2.5} concentration and aerosol optical depth over the
581 background of North China from 2009 to 2011. *Atmos. Res.* 128, 179–188.

582 Yao, F., Si, M., Li, W., Wu, J., 2018a. A multidimensional comparison between MODIS and
583 VIIRS AOD in estimating ground-level PM_{2.5} concentrations over a heavily polluted
584 region in China. *Sci. Total Environ* 618, 819-828.

585 Yao, J., Raffuse, S.M., Brauer, M., Williamson, G.J., Bowman, D.M., Johnston, F.H.,
586 Henderson, S.B., 2018b. Predicting the minimum height of forest fire smoke within the
587 atmosphere using machine learning and data from the CALIPSO satellite. *Remote Sens.*
588 *Environ.* 206, 98–106.

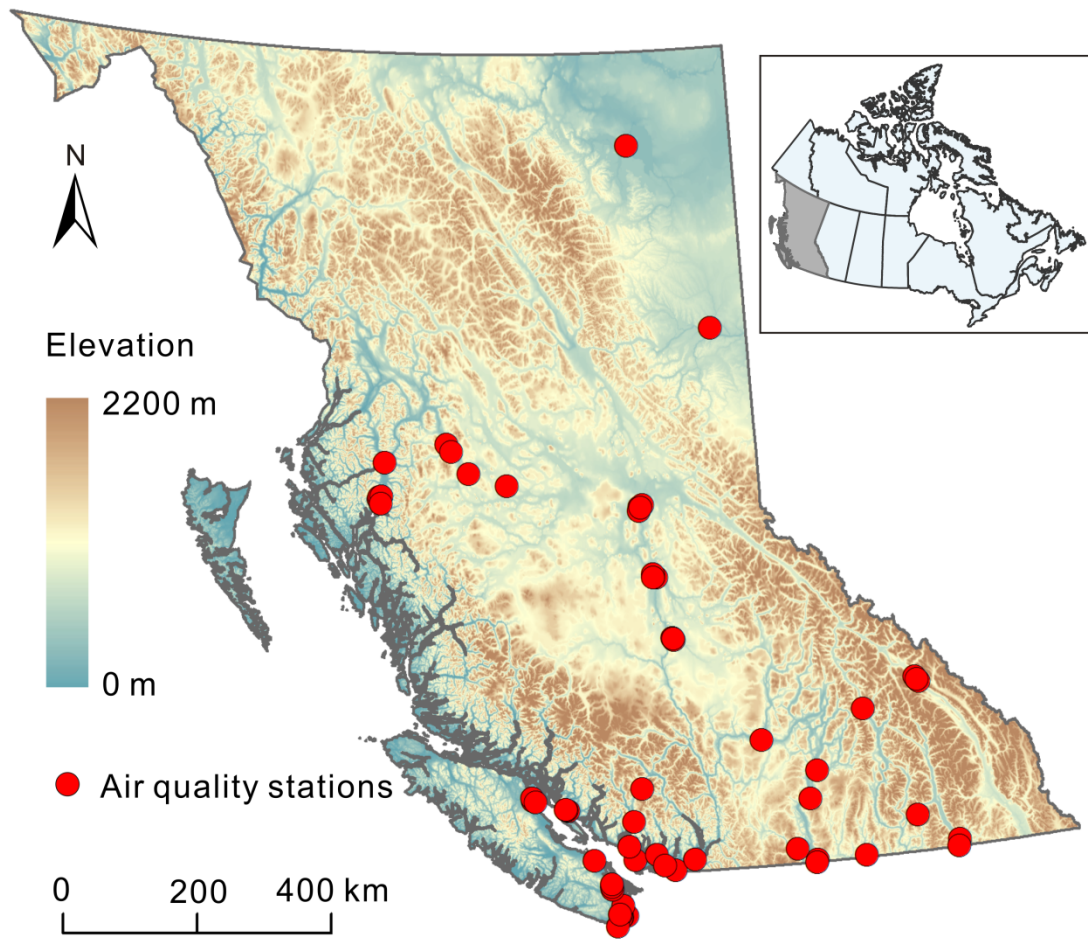
589 You, W., Zang, Z., Pan, X., Zhang, L., Chen, D., 2015. Estimating pm_{2.5} in Xi'an, China using
590 aerosol optical depth: a comparison between the MODIS and MISR retrieval
591 models. *Sci. Total Environ.* 505, 1156–1165.

592 You, W., Zang, Z., Zhang, L., Li, Y., Pan, X., Wang, W., 2016. National-scale estimates of

593 ground-level PM2.5 concentration in China using geographically weighted regression
594 based on 3 km resolution MODIS AOD. *Remote Sens.* 8, 184.

595 Zheng, Y., Zhang, Q., Liu, Y., Geng, G., He, K., 2016. Estimating ground-level PM2.5
596 concentrations over three megalopolises in China using satellite-derived aerosol optical
597 depth measurements. *Atmos. Environ.* 124, 232-242.

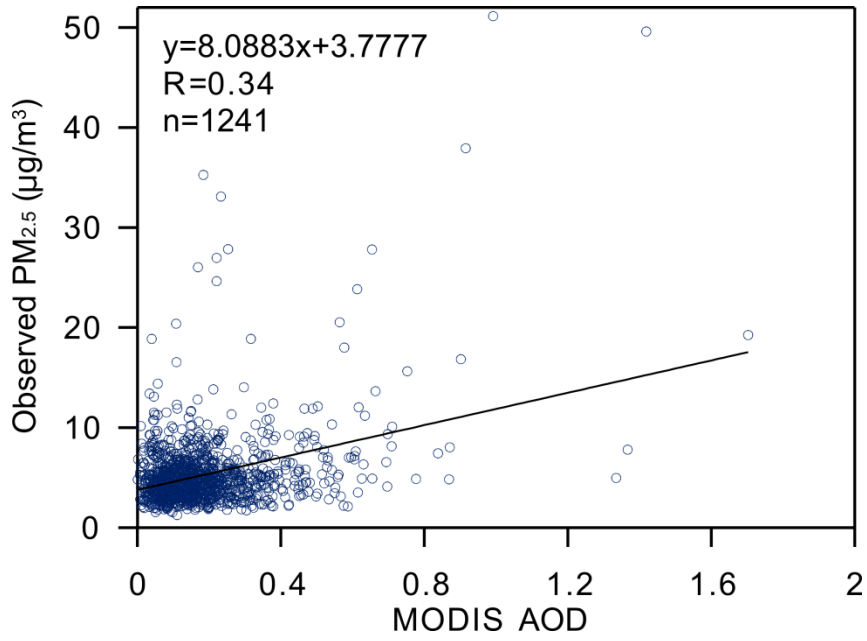
598 Zhan, Y., Luo, Y., Deng, X., Chen, H., Grieneisen, M.L., Shen, X., Zhu, L., Zhang, M., 2017.
599 Spatiotemporal prediction of continuous daily PM2.5, concentrations across China
600 using a spatially explicit machine learning algorithm. *Atmos. Environ.* 155, 129-139.



601

602

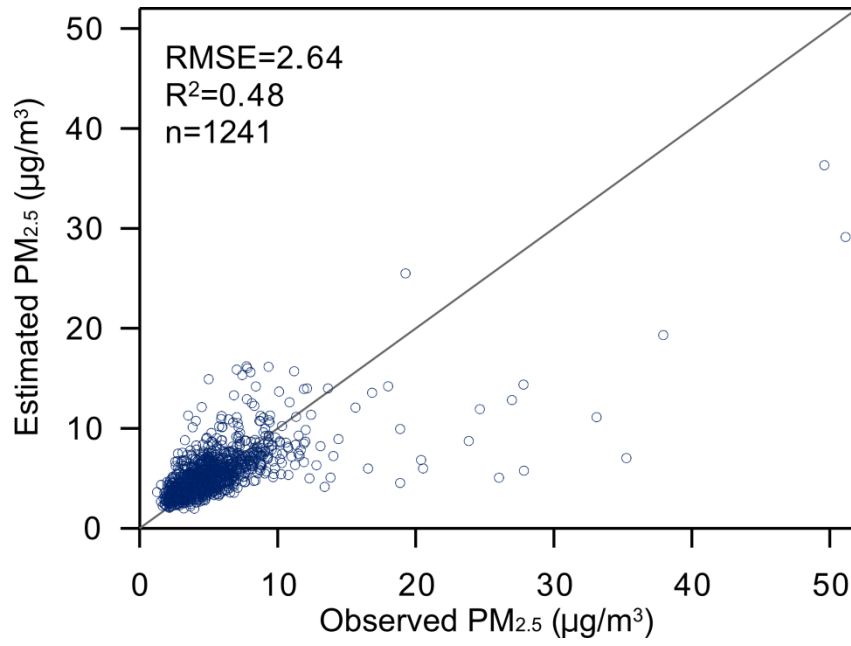
Fig. 1. Study Site. Red dots represent the location of air quality stations across BC.



603

604 Fig.2 Empirical relationship between PM_{2.5} and AOD. X-axis indicated the AOD values derived

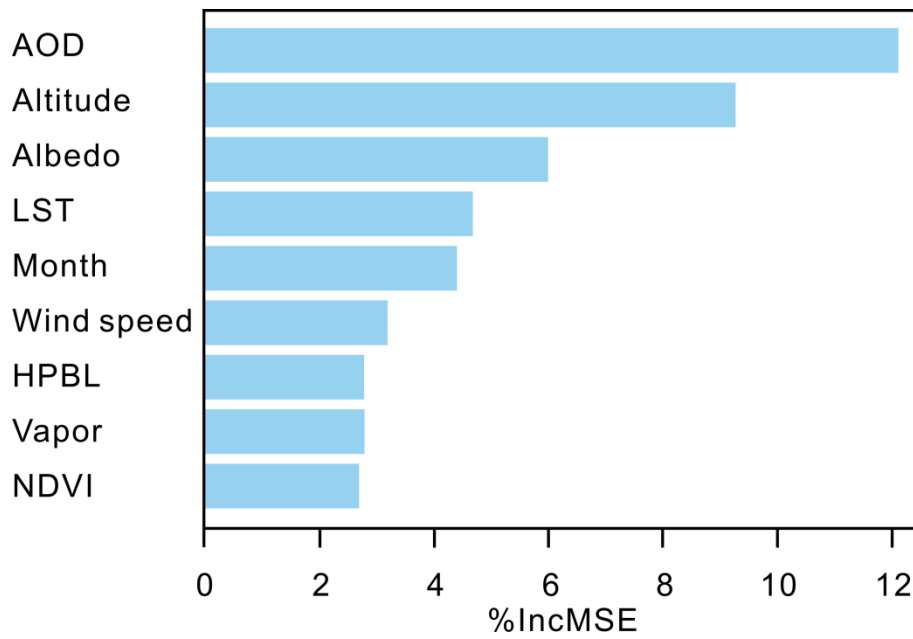
605 from MODIS dataset. Y-axis indicated the PM_{2.5} retrieved from the air quality stations.



606

607

Fig. 3 Comparison between observed and estimated PM_{2.5} using Cubist.



608

609

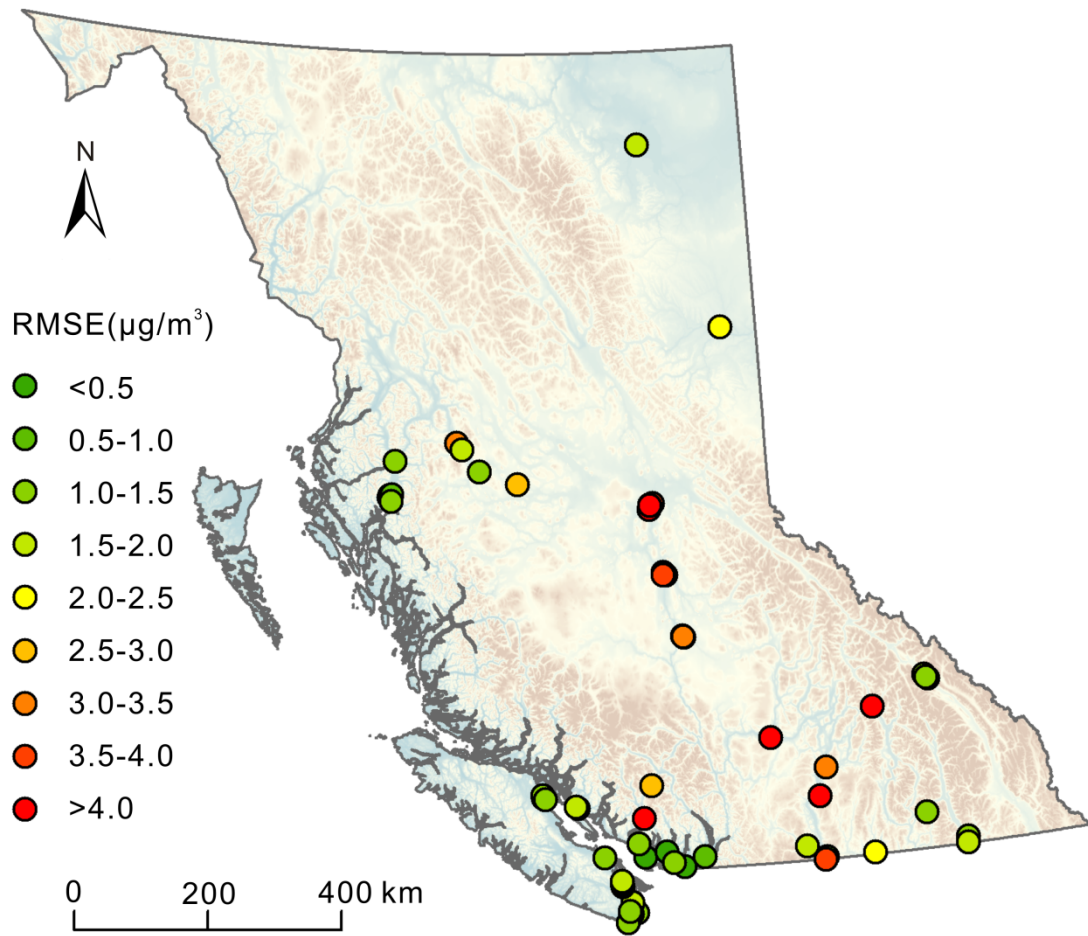
Fig. 4 Variable importance analysis (Cubist Model). Y-axis indicated the predictors for

610

predicting $PM_{2.5}$. X-axis indicated the percentage increase in mean square error (%IncMSE)

611

without using the corresponding predictor.

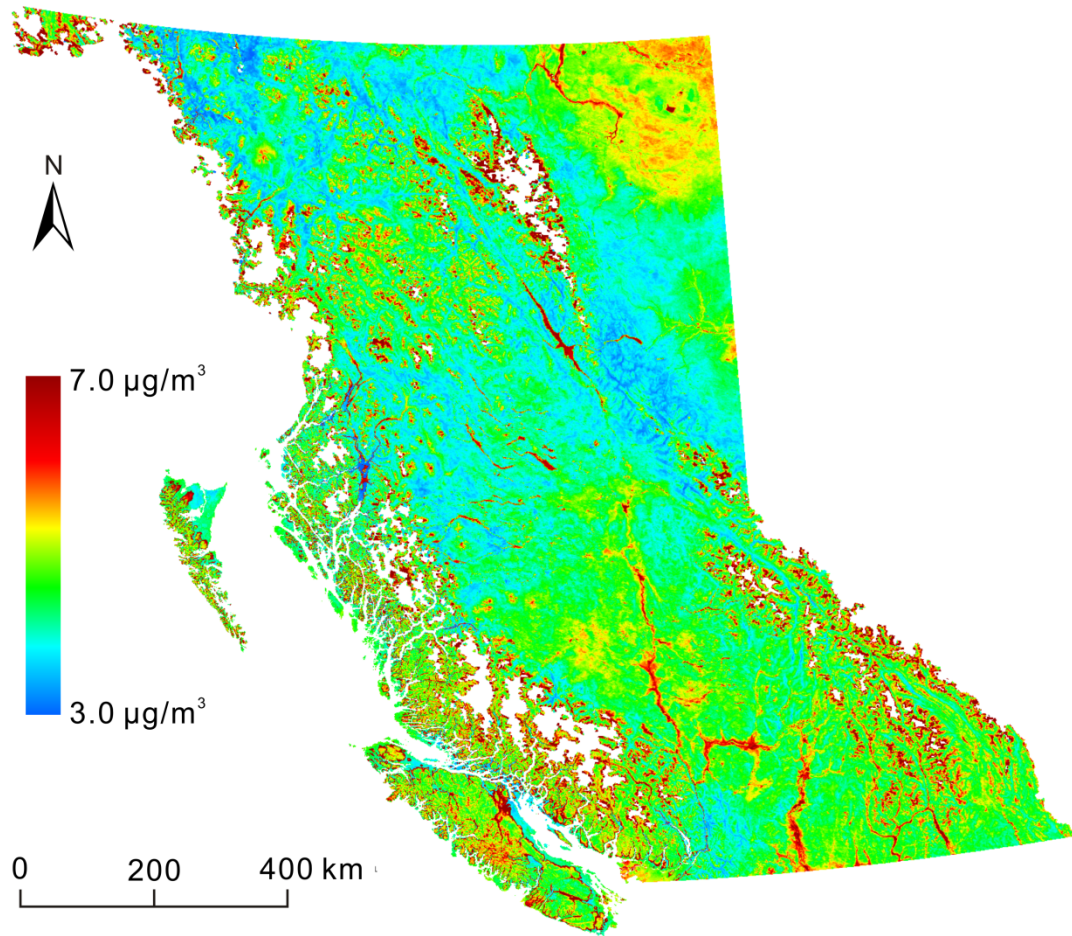


612

613 **Fig. 5** Location-based root mean square error (RMSE) of estimated $\text{PM}_{2.5}$. Red indicated an

614 air quality station with higher RMSE, and green indicated a station with lower RMSE after a

615 comparison with observed data.



616

617

Fig. 6 Average of ground-level PM_{2.5} concentration across BC (2001-2014)

618 **Table 1** Information on datasets used for PM2.5 estimation

Dataset	Spatial resolution	Temporal resolution	Scenes	Derived predictors
MOD04_3k	3km	Daily	25350	AOD
MOD05_L2	1km	Daily	22198	Vapor
MOD11A1	1km	Daily	25369	LST
MOD13A3	1km	Monthly	1677	NDVI
MCD43B3	1km	16 days	6394	albedo
NCAR/NCEP re-analysis	2.5°	Monthly	/	HPBL, wind speed
SRTM DEM	90m	/	/	elevation

619

620 **Table 2** Accuracy of PM_{2.5} prediction of each machine learning model.

Model	CV-RMSE ($\mu\text{g}/\text{m}^3$)	CV-R ²
MLR	3.24	0.22
BRNN	3.04	0.31
SVM	3.13	0.30
LASSO	3.20	0.24
MARS	3.05	0.31
RF	2.67	0.49
XGBoost	2.71	0.46
Cubist	2.64	0.48

621