1

**Estimation of surface PM$_{2.5}$ concentrations from atmospheric gas species retrieved**

**from TROPOMI using deep learning:**

**Impacts of fire on air pollution over Thailand**

Rackhun Son[a,d,e], Hyun Cheol Kim[b,c], Jin-Ho Yoon[d] and Dimitris Stratoulias[e,f,]*

[a]*Department of Biogeochemical Integration, Max Planck Institute for Biogeochemistry, Jena, Germany*

[b]*Air Resources Laboratory, National Oceanic and Atmospheric Administration, College Park, MD, USA*

[c]*Cooperative Institute for Satellite Earth System Studies, University of Maryland, College Park, MD, USA*

[d]*School of Earth Science and Environmental Engineering, Gwangju Institute of Science and Technology, Gwangju, Korea*

[e]*Geospatial Information department, Asian Disaster Preparedness Center (ADPC), Bangkok, Thailand*

[f]*SERVIR-Mekong, Bangkok, Thailand*

* Dimitris Stratoulias: dimitrios.stratoulias@adpc.net

1

**Abstract**

20

21   Surface PM$_{2.5}$ concentration is routinely observed at limited number of surface

22   monitoring stations. To overcome its limited spatial coverage, space-borne

23   monitoring system has been established. However, it also faces various challenges

24   such as cloud contamination and limited vertical resolution. In this study, we propose

25   a deep learning-based surface PM$_{2.5}$ estimation method using the attentive

26   interpretable tabular learning neural network (TabNet) with atmospheric gas species

27   retrieved from the tropospheric monitoring instrument (TROPOMI). Unlike previous

28   applications that primarily used decision tree-based algorithms, TabNet provides

29   interpretable decision-making steps to identify dominant factors. By incorporating

30   five TROPOMI products (i.e., NO$_2$, SO$_2$, O$_3$, CO, HCHO), we have tested the

31   system's capability to produce surface PM$_{2.5}$ concentration without aerosol optical

32   property, which was used more traditionally. The proposed model successfully

33   captures spatiotemporal variations and its performance surpasses those of other

34   leading machine learning models over Thailand in the period of 2018-2020. The

35   interpretable decision-making steps highlight that carbon monoxide is the most

36   influential chemical component, which relates to the seasonal burning in southeast

37   Asia. It is found that air quality impacts from fire are stronger in the northern part of

38   Thailand and fires in neighboring countries should not be neglected. The proposed

39   method successfully estimates surface PM2.5 concentration without aerosol optical

40   property, implying its potential to advance monitoring air quality over remote

41   regions.

42   Keywords: PM$_{2.5}$; TROPOMI; deep learning; TabNet

43

44

2

45    Key Points

46        1.  TabNet successfully estimates surface $PM_{2.5}$ with atmospheric gas compositions

47        2.  CO is highlighted as a key factor in estimating spatiotemporal pattern of $PM_{2.5}$

48        3.  The origin of CO is likely from seasonal fire in Thailand and its bordering countries

49

## 1. Introduction

Ambient air pollution is a critical threat to public health, causing more than three million premature fatalities worldwide in 2012 (Organization, 2016) as well as various environmental issues (Gurjar et al., 2010). Among air pollutants, fine particulate matter ($PM_{2.5}$), namely, ambient airborne particulates sized under 2.5 microns, is well known to damage human health seriously. Due to its microscopic size, $PM_{2.5}$ can affect the respiratory and cardiovascular systems, causing or worsening major illnesses such as asthma, lung cancer and heart disease (Weichenthal et al., 2013). Rising public concern about air quality urges not only reductions in air pollutants, but also improvements to air quality monitoring at the ground level to assess the health and socioeconomic impacts.

Annual mean $PM_{2.5}$ concentrations in Thailand reached 21.4 µg/m³ in 2020, making it the 34th most polluted country in the world ("World Air Quality Report 2020," n.d.). An estimated 40,000 deaths annually in Thailand are attributable to ambient air pollution (Pinichka et al., 2017), resulting in 0.74 to 1.33 million USD worth of economic costs (Vassanadumrongdee and Matsuoka, 2005). Although air pollution exposure in Thailand temporarily improved during the recent COVID-19 pandemic (Rodríguez-Urrego and Rodríguez-Urrego, 2020; Stratoulias and Nuthammachot, 2020), it remains high due to widespread smoke emissions from agricultural burnings and forest fires (Punsompong et al., 2021). $PM_{2.5}$ emissions from burning crop residue and forest fires are estimated to be 141,000 and 5,000 tons per year, respectively, mostly concentrated in the central and northern regions of Thailand (Junpen et al., 2013; Kanabkaew and Kim Oanh, 2011). However, the insufficient number of in-situ $PM_{2.5}$ measurements, especially for the provinces in the north and northeast of the country (Figure 1a), limits monitoring air quality

73    and establishing a national plan for its management. Given that considerable financial and

74    time resources are required to increase the number of air quality monitoring stations,

75    satellite remote sensing-based $PM_{2.5}$ estimation is an alternative way to increase the limited

76    spatial coverage.

77          Two different methods have been developed and widely used to estimate surface

78    $PM_{2.5}$ concentration: Chemical Transport Models (CTMs) and statistical regression models.

79    Based on physicochemical processes and atmospheric conditions, chemical transport

80    models can approximate the quantity of air pollutants with continuous spatiotemporal

81    coverage (Liu et al., 2004; Van Donkelaar et al., 2010). However, uncertainties in emission

82    inventories and limited representation of chemical reactions in the ambient atmosphere

83    remain major concerns (Shin et al., 2020). Among statistical approaches, multiple linear

84    regression has been the most commonly applied in the early stages (Chu et al., 2016). Also,

85    geographically weighted regression, an extension of multiple linear regression, has been

86    proposed to assign distance-based weights to reflect spatial variability and local effects to

87    provide regional estimations (Brunsdon et al., 1998; Jiang et al., 2017; You et al., 2016).

88    Mixed-effect models adopt fixed and random effect terms to separate statistical relationship

89    and variability by time and region (Kloog et al., 2012; Xie et al., 2015). In addition, the

90    generalized additive model has been proposed to consider the nonlinear characteristics

91    between input and target variables (Sorek-Hamer et al., 2013; Zou et al., 2017).

92          Machine learning (ML) algorithms have recently introduced as innovative

93    developments in the bottom-up approaches to upscale data-driven in-situ models to

94    spatially explicit gridded estimates. Random forest (RF), one of the most frequently applied

95    algorithms, has further improved estimation accuracy and has higher interpretability at both

96    national and regional scales (Chen et al., 2018; Hu et al., 2017; Wei et al., 2019). Elastic-

97  net application has successfully expanded the spatiotemporal dimension with a large

98  number of predictors (Xue et al., 2019). Support vector machine (SVM) can enhance spatial

99  resolution at a 100 m scale by being merged into multiple modeling stages (de Hoogh et al.,

100 2018). Other ML models such as Bayesian maximum entropy (Jiang and Christakos, 2018),

101 gradient boosting machine models (Chen et al., 2019; Wang et al., 2021) and RF combined

102 with ordinary kriging (Han et al., 2022) have also been employed to incorporate satellite-

103 derived products into ground-level observations.

104     As computing technology and resources have advanced, neural network-based

105 approach has introduced deeper and wider layers, defined as deep learning (DL), and has

106 begun to outperform classical ML models based on decision tree algorithms in various

107 regression tasks (Devlin et al., 2018; He et al., 2016). DL based methods have also recently

108 been attempted in remote sensing due to their high accuracy using large amounts of data

109 (Ghahremanloo et al., 2021; Zhang et al., 2020; Zhu et al., 2020). However, compared with

110 decision trees, the usability of this cutting-edge approach is yet to be explored in-depth for

111 $PM_{2.5}$ satellite-based estimation.

112     Thus, this study aimed to develop a DL-based model to estimate daily ground-level

113 $PM_{2.5}$ concentrations based in-situ observations in Thailand and satellite-derived

114 atmospheric gas products. Regarding the DL network architecture, we implemented the

115 Attentive Interpretable Tabular Learning neural network (TabNet) (Arık and Pfister, 2021),

116 which is tailored for use with tabular datasets. We evaluated the model's performance

117 through five different regions in Thailand and compared it with other popular machine

118 learning algorithms such as SVM, RF, XGBoost (Chen and Guestrin, 2016), LightGBM

119 (Ke et al., 2017) and CatBoost (Prokhorenkova et al., 2018). Furthermore, to shed light on

120   the critical characteristics of $PM_{2.5}$ concentration in Thailand, we analyzed the global/local

121   feature selection and reasoning processes as well as fire impacts on $PM_{2.5}$ concentration.

122

123   **2. Study area and Data**

124   *2.1. Study area*

125        Thailand is located at the center of the Indochinese peninsula, has the 10th largest

126   economy in Asia and hosts a population of almost 70 million people ("World Economic

127   Outlook (October 2021)," n.d.). The country is divided into 76 administrative provinces as

128   primary local government units and the capital Bangkok. In this study, Thailand is divided

129   into five regions to analyze the country's regional characteristics: north, northeast, central,

130   east and south (Figure 1a). Despite the low air quality in Thailand, ground monitoring

131   stations are sparse and mostly concentrated in the central region, which contains 29 of the

132   67 stations used in this study (green points in Figure 1a). For the remaining regions, 15, 5,

133   11 and 7 stations are distributed in the north (blue), northeast (red), east (yellow) and south

134   (magenta), respectively.

135   *2.2. Ground-level $PM_{2.5}$ observation*

136        The Pollution Control Department in the Air Quality and Noise Management

137   Bureau provides national air quality monitoring records for approximately 84 stations (as of

138   2021). Considering the consistency of the data availability during the experimental period

139   from January 2018 to June 2021, we selected the daily measurements of $PM_{2.5}$

140   concentration from 67 stations (Figure 1a) as the target dataset for the model training. The

141   observed $PM_{2.5}$ concentration pattern has an exponential quantile-quantile distribution

7

142　(Figure 2a). This asymmetry can hamper model training by blurring the variance in the

143　pollution levels over different input conditions. To transform the data to be closely fitted by

144　a normal distribution, we thus took logs of the $PM_{2.5}$ values after adding one (Figure 2b)

145　and the results showed significantly higher R-squared coefficients ($R^2$) from 0.744 to 0.996.

146　***2.3. TROPOMI***

147　　　　The Sentinel-5P mission is a precursor satellite measuring atmospheric chemical

148　concentrations at high spatial and radiometric resolutions. The TROPOspheric Monitoring

149　Instrument (TROPOMI) onboard Sentinel-5P is designed to record the reflectance of

150　wavelengths using multispectral sensors. We utilized five TROPOMI products (Borsdorff

151　et al., 2018; De Smedt et al., 2018; Garane et al., 2019; Theys et al., 2017; Van Geffen et

152　al., 2019): the tropospheric $NO_2$ column ($NO_2$), $SO_2$ vertical column density at the ground

153　level ($SO_2$), total atmospheric column of $O_3$ ($O_3$), vertically integrated column of CO (CO)

154　and tropospheric formaldehyde column (HCHO); this was based on 354 of the 388

155　wavelength pairs. TROPOMI Level 2 products are accessible from the Copernicus Open

156　Access Hub website (https://s5phub.copernicus.eu), and we retrieved a daily Level 3 pre-

157　processed dataset from the Google Earth Engine using the quality assurance values of 0.75

158　for $NO_2$ and 0.5 for the other components except for $O_3$ and $SO_2$. The Sentinel-5P images

159　were co-located with the ground station data and the values of the pixel encompassing the

160　point location of the ground station were extracted to train the model. When the spatial

161　mapping of $PM_{2.5}$ was inferred, the datasets were resampled to a 10 km grid to incorporate

162　other auxiliary datasets. Subsequently, the variables, except for $O_3$, were transformed into a

163　logarithmic scale similar to $PM_{2.5}$. Considering that the ranges of each variable varied,

164　specified constants were multiplied and added before the log transform (Figures S1a–d).

165 *2.4. Meteorological dataset*

166      ERA5-Land (Muñoz Sabater, 2019) provides a dataset for land components from

167 ERA5, the fifth-generation climate reanalysis dataset provided by the Copernicus Climate

168 Change Service at the European Centre for Medium-Range Weather Forecasts. Following

169 previous studies (Chen et al., 2018; Wei et al., 2019), we adopted seven meteorological

170 components from the reanalysis dataset: temperature and dew-point temperature at a 2 m

171 height, total evaporation, surface pressure, precipitation and wind components at a 10 m

172 height. We also approximated relative humidity and wind speed using Eqs. (1) and (2):

173
$$relhumidity = 100 \times \frac{e^{\frac{17.625 \times T_d}{243.04 + T_d}}}{e^{\frac{17.625 \times T}{243.04 + T}}} \tag{1}$$

174
$$windspeed = \sqrt{U^2 + V^2} \tag{2}$$

175 where $T$ is temperature, $T_d$ is dew-point temperature, $U$ is the horizontal wind component

176 (U-wind) and $V$ is the meridional wind component (V-wind). For precipitation and wind

177 speed, the scaled log transform was applied as mentioned above (Figures S1e, and f).

178 Furthermore, we considered geographical factors such as elevation from ETOPO1 (Amante

179 and Eakins, 2009) with a 1 arc-minute resolution to integrate the land topography and

180 bathymetry and land cover classifications from GlobCover (Arino, 2010). These were

181 categorized into 22 types based on observations from the ENVISAT satellite mission for

182 2009 with a spatial resolution of approximately 300 m.

183

## 3. Methodology

### 3.1. TabNet

TabNet is a novel neural network architecture designed to provide an adequate tabular dataset (Arık and Pfister, 2021). Based on an encoder/decoder structure, high-dimensional features can be transformed into a meaningful representation through trainable embedding layers without any pre-processing steps. For instance, the layers can map categorical features into a numerical format as well as handle raw numerical features without normalizing global features. One salient strategy of the TabNet is to employ the sequential attentive transformer architecture to select the importance features in decision steps. In each step, learnable masks search for a subset of the relevant features by quantifying the contribution of the decision.

### 3.2. Interpretability

The feature attribution mask $\mathbf{M} \in \mathbb{R}^{B \times D}$ provides instance-wise interpretable insights for reasoning; $B$ is the batch size and $D$ is the dimension of the feature. At the $i^{th}$ decision step, the processed features from the preceding step $\mathbf{a}[i\text{-}1]$ are given to a trainable nonlinear processing $h_i$, composed of a fully connected layer, batch normalization and gated linear unit (Dauphin et al., 2017). The mask is obtained through a sparse regulation function, which we set using *entmax* (Peters et al., 2019), as summarized in Eq. (3):

$$M[i] = entmax(P[i-1] * h_i(a[i-1]))$$  (3)

$\mathbf{P}[i]$ is the prior scale term to regulate the flexibility of feature selection in the multiple steps, as defined in Eq. (4):

205
$$P[i] = \prod_{j=1}^{i}(\gamma - M[j]) \qquad (4)$$

206 where $\gamma$ is the coefficient for the feature reselection in the mask. $\mathbf{P}[0]$ is initialized as all

207 ones, $\mathbf{1}^{B \times D}$, indicating that none of the features are used at the beginning. As a feature is

208 considered thoroughly, its scale term is reduced to focus on the other features in the next

209 steps. The weights of the trained mask represent the relative importance of each step in all

210 instances. For example, if $\mathbf{M}_{b,j}[i]$=0, then the $j^{th}$ feature should have no decision

211 contribution in the $i^{th}$ step for the $b^{th}$ sample. Finally, the aggregated weights from the

212 masks allow us to understand the importance of each feature in terms of its global behavior.

213 ### 3.3. Training details

214     The weather in Thailand has distinct seasonality; the rainy season, which usually

215 lasts from June to October, can significantly affect the $PM_{2.5}$ concentration in the

216 atmosphere (Figure S2). Moreover, the mapping of averaged $PM_{2.5}$ displays a higher

217 concentration in the northern area, above 40, than elsewhere (Figure 1b). Considering these

218 spatiotemporal characteristics, we added the observed month and geographical coordinates

219 (longitude and latitude) of the station as input features. In total, 19 input variables were

220 used in this study: $NO_2$, $SO_2$, $O_3$, CO, HCHO, temperature, dew-point temperature, relative

221 humidity, U-wind, V-wind, wind speed, precipitation, pressure, evaporation, elevation, land

222 cover type, month, longitude and latitude. For the categorical variables, namely, month and

223 land cover type, we set the embedding dimensions to 6 and 17, respectively.

224     Following convention, we randomly split the data from 2018 to 2020 into training

225 and testing datasets using an 80:20 ratio; the number of samples were 14,069 and 3518,

226 respectively. We also evaluated the functionality of upscaled mapping using a 10 km

227  resolution grid format of the input dataset for 2021. To ensure robust training, a 5-fold

228  cross-validation was set, and the final $PM_{2.5}$ estimation was calculated by averaging the

229  results from the five trained models. The model was implemented using the *pytorch_tabnet*

230  package (https://github.com/dreamquark-ai/tabnet) and trained with the Adam algorithm

231  with weight decay using a 0.01 learning rate and a batch size of 64. Following the

232  guidelines for hyperparameters (Arık and Pfister, 2021), we set the depth and width of

233  TabNet as follows: $N_d = N_a = 24$, $N_{steps} = 4$, $\gamma = 1.3$ and $\lambda_{sparse} = 0.001$.

234

235  **4. Results**

236  *4.1. Evaluation of general model performance*

237  Figure 3 presents the accuracy validation results of the estimated $PM_{2.5}$

238  concentration for Thailand and the five divided regions. For the entire study domain (Figure

239  3a), three evaluation metrics show 0.873 of $R^2$, 9.22 of root mean square error (RMSE) and

240  20.62 of the mean absolute percentage error (MAPE). When these results are compared

241  with other state-of-the-art ML algorithms, $R^2$ and RMSE of the proposed method show the

242  best scores (Table 1). In terms of the linear relationship between the observations and

243  estimated $PM_{2.5}$ concentrations, all the models show slope coefficient values under 1. These

244  results imply that the ML models tend to underestimate the $PM_{2.5}$ concentration, as is

245  consistently reported in previous studies (Ma et al., 2016; Wei et al., 2019). TabNet can

246  compensate for this bias, as it has the highest value of the slope coefficient (0.84). This

247  improvement is especially noticeable in the extremely high concentration cases of more

248  than 300 µg/m³ (Figure S3).

12

249        When the scores of evaluation metrics are compared by region, the highest value of

250    $R^2$ (0.884) is observed in the north (Figure 3b). These results are consistent with the

251    mapping of $R^2$ for each station showing higher than 0.8 of $R^2$ in all the stations in the north,

252    including the Chiang Mai and Lampang provinces (Figure S4a). On the other hand, the

253    scale of biases is larger than other regions with 13.44 of RMSE, due to its wider range of

254    the $PM_{2.5}$ concentration exceeding 300 µg/m$^3$ as a maximum (Figures 3b and S4b). Given

255    that the $PM_{2.5}$ concentrations in the north are generally higher (Figure 1b) and extreme

256    cases are more frequent due to agricultural burnings and forest fires (Punsompong et al.,

257    2021), the large errors are typically caused by the underestimation mentioned previously,

258    particularly for high concentration cases. When the regional differences in scale are

259    diminished by considering the ration of the scale between the errors and actual values,

260    some stations in Bangkok and neighbor cities show higher scores of MAPE (Figure S4c).

261    But the south region shows the lowest accuracy with 21.51 of MAPE and 0.507 of $R^2$

262    (Figure 3f). The distinctively low slope coefficient in the south represents that its poor

263    performance is mainly caused by underestimation (Figure S4d). Considering that the air

264    quality of southern Thailand is influenced by pollutants from peatland fires in Indonesia

265    during the southwest monsoon (Mahasakpan et al., 2023), our model seems to have limits

266    to estimate air mass transportation from out of the study domain.

267    *4.2. Application on high-coverage mapping*

268        One of the main purposes of employing remote sensing data is to enlarge the spatial

269    coverage of $PM_{2.5}$ monitoring. Figure 4 illustrates the monthly averaged results of the $PM_{2.5}$

270    estimation for 2021. The mapping results (Figures 4a–f) generally agree with the

271    observations (Figures 4g–l) with respect to seasonal variation by region. In January, the

272     central region of Thailand shows high levels of $PM_{2.5}$ concentrations. In the north, the

273     concentrations significantly increase from January and peak at over 60 μg/m$^3$ in March.

274     The regional time difference in the peak of air pollutants can be explained by the fact that

275     harvesting and residue burning are carried out in a different season in each region

276     (Kanabkaew and Kim Oanh, 2011). Thereafter, the concentrations decrease in all the

277     regions as the rainy season approaches.

278        To evaluate the temporal variation of the $PM_{2.5}$, we compare the daily variations in

279     the observed and estimated $PM_{2.5}$ over the five regions of Thailand (Figure 5). The northern

280     area shows the highest performance scores (0.83, 12.58 and 19.81 for $R^2$, RMSE and

281     MAPE, respectively). The value of slope coefficient is almost 1 representing a significant

282     improvement in the underestimation for extreme levels of $PM_{2.5}$, with high accuracy for

283     peak days during March and April. The other regions, except for the south, also show good

284     performances according to the evaluation metrics. Although the south region has smaller

285     scale of error (4.69 of RMSE), the underestimation on high concentration days, particularly

286     on those days with values above 60 μg/m$^3$, has scope for further exploration and

287     improvement for long-range transport effects across neighboring countries.

288

289     **5. Discussion**

290     *5.1. Model interpretation*

291        Interpretability makes it possible for us to understand model's behavior at each

292     learning steps and to point out important processes, which can be translated into more

293     practical way. However, there isn't a perfect method to interpret ML and DL approaches,

294     which is well recognized as a potential limitation. A major advantage of the TabNet is its

14

295     attentive transformer structure, which provides post-hoc explanations by assessing the

296     contribution of each feature from both global and local perspectives. First, the global

297     importance of each feature is illustrated in Figure 6. The observed month displays the

298     highest ratio of contribution with approximately 40% of importance, which is expected

299     according to the seasonality of $PM_{2.5}$ in Thailand (see Figure S2). Geographical features

300     such as land cover type, coordination and elevation follow next, demonstrating their

301     importance. In terms of chemical components, $NO_2$ and $SO_2$, which are commonly known

302     as precursors in the secondary formation of $PM_{2.5}$ (Baker and Scheff, 2007; Tucker, 2000),

303     rank relatively low among all the features; $SO_2$ shows almost zero contribution to the

304     estimation. Instead, CO accounts for about 20% of the contribution. Considering that CO is

305     a by-product of carbon-containing fuel combustion, these results agree with the scenario

306     that vehicular emissions and fires have a greater impact on the variation in air quality in

307     Thailand than industrial emissions (ChooChuay et al., 2020).

308         Figure S5 illustrates the top five important features on each decision step as the

309     aspect of local feature importance. Consistent with the global perspective, the observed

310     month, CO and land cover type are ranked as the most determining factors in all the steps,

311     regardless of season. Interestingly, the second step displays different composition of

312     importance, especially for meteorological features, by season. The importance of wind

313     speed and relative humidity are relatively lower for dry season ranking fourth and fifth

314     (Figure S5f), while they are selected as the second and third most important features in wet

315     season (Figure S5j). Some other meteorological factors, such as pressure, evaporation and

316     dew-point temperature, are also displayed in other steps, in spite of their low contribution

317     (less than 5%). Considering that windy and humid weather can reduce pollution levels, the

318     trained model locally employs weather information to identify the ideal conditions for

319     lower $PM_{2.5}$ concentrations.

### 5.2. Impacts of fire on PM2.5 concentration in Thailand

321         To investigate the impact of fire on the air quality in Thailand, we analyze spatial

322     distribution of fire radiative power (FRP) from the Global Fire Assimilation System

323     (GFAS) in the Copernicus Atmosphere Monitoring Service (CAMS) and chemical

324     components (Figure 7) for a period when all the sub-regions show the rise of PM2.5

325     concentration (from February 25th to March 2nd 2021, red columns in Figure 5). During this

326     period, high levels of FRP were mainly observed in the central region and the border areas

327     in the north and east adjacent to neighboring countries, such as Myanmar, Laos and

328     Cambodia. The concentrations of $PM_{2.5}$ and major chemical components also increased

329     nearby fire hotspots, especially when the highest FRP was observed in the central west of

330     Thailand on March 1st. This causal relationship between FRP and the concentrations can be

331     seen in all the regions during the dry season for the year 2021. FRP shows statistically

332     significant correlations with $PM_{2.5}$ in the north, northeast and east regions as well as with

333     other chemical components in the north (Figures S6-10). Considering that those regions

334     generally have higher concentrations of $PM_{2.5}$, the results demonstrate that the frequency

335     and duration of fire can significantly influence on the level of air quality in Thailand.

336     Besides, fires in the neighbor countries can also be another factor to cause considerable

337     increase of $PM_{2.5}$ concentration. For instance, when many hotspots were detected in the

338     territory of Myanmar and Laos on February 26th and Marth 1st, the chemical and $PM_{2.5}$

339     concentrations distinctly increased in the north and northeast parts of Thailand the next day.

16

### 5.3. Potential further improvements

Traditionally, aerosol optical depth (AOD) has been played as an essential factor to estimate the surface level of $PM_{2.5}$ concentrations. However, the presence of cloud and snow along with the limited vertical resolution causes unfeasibility for archive reliable AOD limiting its spatial coverage (Hsu et al., 2013; Levy et al., 2007). Our approach based on atmospheric gas composition offers a viable alternative to address the spatial limits. We also tested the skill of estimation including aerosol index and its results do not show any clear difference in the metric scores (Figure S11), proving that $PM_{2.5}$ concentration can be accurately estimated only with atmospheric trace gases. Although the spatial constraints still remain in this study due to excluding data below the quality threshold, future work will focus on the development of a model to handle the low-quality data aiming to achieve reliable full coverage for $PM_{2.5}$ estimation.

While there have been prior attempts to apply DL-based modeling to estimate $PM_{2.5}$, its performance is lower than that of other algorithms (Chen et al., 2022; Pu and Yoo, 2021; Wong et al., 2021). Importantly, these results could be linked to their simple model structures, which mostly consist of a series of fully connected hidden layers with nonlinear activation functions. In the current study, we showcase how by fusing TROPOMI data with other geospatial sources and incorporating an advanced DL algorithm to provide an accurate representation of $PM_{2.5}$ concentration; consequently, an air pollution indicator can be developed. Previous studies have reported the potential of DL algorithms such as CNN and LSTM to improve estimation performance (Chen et al., 2021; Lu et al., 2021), and our results also support this by adopting a state-of-the-art DL algorithm. Numerous advanced DL methods have recently been developed and have achieved remarkable progress in several fields (Devlin et al., 2018; He et al., 2016); however, applying DL to

17

364　estimate $PM_{2.5}$ concentration has not yet been widely explored. Thus, monitoring air quality

365　by implementing DL approaches has considerable room for improvement.


366

**6. Conclusion**

368　　　To estimate ground-level $PM_{2.5}$ concentration across Thailand, we develop a novel

369　method based on DL algorithm, namely TabNet, with atmospheric gas from the TROPOMI

370　on the Sentinel-5P. Our model shows more robust performance than other state-of-the-art

371　ML algorithms, with an $R^2$ and RMSE (MAPE and slope coefficient) of 0.873 and 9.22

372　(20.62 and 0.84), respectively. The interpretable decision processes in TabNet indicate that

373　monthly variation is the most significant feature in $PM_{2.5}$ estimation. Geospatial

374　characteristics such as land cover type and latitude also provide a notable contribution from

375　a global perspective. Among the chemical components from the TROPOMI, CO shows

376　higher ratio of importance than the others. These results suggest that emissions from

377　biomass burning influence air quality in Thailand considerably. In the low-level $PM_{2.5}$

378　concentration scenarios, humid and windy weather conditions are also highlighted in the

379　local decision processes. Based on its robust performance, the model is applied to generate

380　grid format mapping of $PM_{2.5}$ concentrations. We find that it can capture the temporal

381　variation in and uneven spatial distribution of $PM_{2.5}$ concentrations using a 10 km grid. The

382　enhanced estimation ability and its application are expected to not only boost other air

383　quality studies, but also contribute to air quality management by providing advanced

384　monitoring and evaluation techniques.

385 **Acknowledgements**

388
389

390  **References**

391  Amante, C., Eakins, B.W., 2009. ETOPO1 arc-minute global relief model: procedures, data

392      sources and analysis.

393  Arino, O., 2010. GlobCover 2009.

394  Arık, S.O., Pfister, T., 2021. Tabnet: Attentive interpretable tabular learning, in: AAAI. pp.

395      6679–6687.

396  Baker, K., Scheff, P., 2007. Photochemical model performance for PM2. 5 sulfate, nitrate,

397      ammonium, and precursor species SO2, HNO3, and NH3 at background monitor

398      locations in the central and eastern United States. Atmos. Environ. 41, 6185–6195.

399  Borsdorff, T., Aan de Brugh, J., Hu, H., Aben, I., Hasekamp, O., Landgraf, J., 2018.

400      Measuring carbon monoxide with TROPOMI: First results and a comparison with

401      ECMWF-IFS analysis data. Geophys. Res. Lett. 45, 2826–2832.

402  Brunsdon, C., Fotheringham, S., Charlton, M., 1998. Geographically weighted regression.

403      J. R. Stat. Soc. Ser. D (The Stat. 47, 431–443.

404  Chen, B., Song, Z., Pan, F., Huang, Y., 2022. Obtaining vertical distribution of PM2. 5

405      from CALIOP data and machine learning algorithms. Sci. Total Environ. 805, 150338.

406  Chen, B., You, S., Ye, Y., Fu, Y., Ye, Z., Deng, J., Wang, K., Hong, Y., 2021. An

407      interpretable self-adaptive deep neural network for estimating daily spatially-

408      continuous PM2. 5 concentrations across China. Sci. Total Environ. 768, 144724.

409  Chen, G., Wang, Y., Li, S., Cao, W., Ren, H., Knibbs, L.D., Abramson, M.J., Guo, Y.,

410      2018. Spatiotemporal patterns of PM10 concentrations over China during 2005–2016:

411      a satellite-based estimation using the random forests approach. Environ. Pollut. 242,

412      605–613.

413  Chen, T., Guestrin, C., 2016. Xgboost: A scalable tree boosting system, in: Proceedings of

414      the 22nd Acm Sigkdd International Conference on Knowledge Discovery and Data

415      Mining. pp. 785–794.

416    Chen, Z.-Y., Zhang, T.-H., Zhang, R., Zhu, Z.-M., Yang, J., Chen, P.-Y., Ou, C.-Q., Guo,

417      Y., 2019. Extreme gradient boosting model to estimate PM2. 5 concentrations with

418      missing-filled satellite data in China. Atmos. Environ. 202, 180–189.

419    ChooChuay, C., Pongpiachan, S., Tipmanee, D., Suttinun, O., Deelaman, W., Wang, Q.,

420      Xing, L., Li, G., Han, Y., Palakun, J., 2020. Impacts of PM2. 5 sources on variations

421      in particulate chemical compounds in ambient air of Bangkok, Thailand. Atmos.

422      Pollut. Res. 11, 1657–1667.

423    Chu, Y., Liu, Y., Li, X., Liu, Z., Lu, H., Lu, Y., Mao, Z., Chen, X., Li, N., Ren, M., 2016.

424      A review on predicting ground PM2. 5 concentration using satellite aerosol optical

425      depth. Atmosphere (Basel). 7, 129.

426    Dauphin, Y.N., Fan, A., Auli, M., Grangier, D., 2017. Language modeling with gated

427      convolutional networks, in: International Conference on Machine Learning. PMLR,

428      pp. 933–941.

429    de Hoogh, K., Héritier, H., Stafoggia, M., Künzli, N., Kloog, I., 2018. Modelling daily

430      PM2. 5 concentrations at high spatio-temporal resolution across Switzerland. Environ.

431      Pollut. 233, 1147–1154.

432    De Smedt, I., Theys, N., Yu, H., Danckaert, T., Lerot, C., Compernolle, S., Van

433      Roozendael, M., Richter, A., Hilboll, A., Peters, E., 2018. Algorithm theoretical

434      baseline for formaldehyde retrievals from S5P TROPOMI and from the QA4ECV

435      project. Atmos. Meas. Tech. 11, 2395–2426.

436    Devlin, J., Chang, M.-W., Lee, K., Toutanova, K., 2018. Bert: Pre-training of deep

437      bidirectional transformers for language understanding. arXiv Prepr. arXiv1810.04805.

438      Garane, K., Koukouli, M.-E., Verhoelst, T., Lerot, C., Heue, K.-P., Fioletov, V., Balis, D.,

439          Bais, A., Bazureau, A., Dehn, A., 2019. TROPOMI/S5P total ozone column data:

440          global ground-based validation and consistency with other satellite missions. Atmos.

441          Meas. Tech. 12, 5263–5287.

442      Ghahremanloo, M., Lops, Y., Choi, Y., Yeganeh, B., 2021. Deep Learning Estimation of

443          Daily Ground-Level NO2 Concentrations From Remote Sensing Data. J. Geophys.

444          Res. Atmos. 126, e2021JD034925.

445      Gurjar, B.R., Molina, L.T., Ojha, C.S.P., 2010. Air pollution: health and environmental

446          impacts. CRC press.

447      Han, S., Kundhikanjana, W., Towashiraporn, P., Stratoulias, D., 2022. Interpolation-Based

448          Fusion of Sentinel-5P, SRTM, and Regulatory-Grade Ground Stations Data for

449          Producing Spatially Continuous Maps of PM2. 5 Concentrations Nationwide over

450          Thailand. Atmosphere (Basel). 13, 161.

451      He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in:

452          Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp.

453          770–778.

454      Hsu, N.C., Jeong, M., Bettenhausen, C., Sayer, A.M., Hansell, R., Seftor, C.S., Huang, J.,

455          Tsay, S., 2013. Enhanced Deep Blue aerosol retrieval algorithm: The second

456          generation. J. Geophys. Res. Atmos. 118, 9296–9315.

457      Hu, X., Belle, J.H., Meng, X., Wildani, A., Waller, L.A., Strickland, M.J., Liu, Y., 2017.

458          Estimating PM2. 5 concentrations in the conterminous United States using the random

459          forest approach. Environ. Sci. Technol. 51, 6936–6944.

460      Jiang, M., Sun, W., Yang, G., Zhang, D., 2017. Modelling seasonal GWR of daily PM2. 5

461          with proper auxiliary variables for the Yangtze River Delta. Remote Sens. 9, 346.

462    Jiang, Q., Christakos, G., 2018. Space-time mapping of ground-level PM2. 5 and NO2

463        concentrations in heavily polluted northern China during winter using the Bayesian

464        maximum entropy technique with satellite data. Air Qual. Atmos. Heal. 11, 23–33.

465    Junpen, A., Garivait, S., Bonnet, S., 2013. Estimating emissions from forest fires in

466        Thailand using MODIS active fire product and country specific data. Asia-Pacific J.

467        Atmos. Sci. 49, 389–400.

468    Kanabkaew, T., Kim Oanh, N.T., 2011. Development of spatial and temporal emission

469        inventory for crop residue field burning. Environ. Model. Assess. 16, 453–464.

470    Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., Liu, T.-Y., 2017.

471        Lightgbm: A highly efficient gradient boosting decision tree. Adv. Neural Inf. Process.

472        Syst. 30.

473    Kloog, I., Nordio, F., Coull, B.A., Schwartz, J., 2012. Incorporating local land use

474        regression and satellite aerosol optical depth in a hybrid model of spatiotemporal PM2.

475        5 exposures in the Mid-Atlantic states. Environ. Sci. Technol. 46, 11913–11921.

476    Levy, R.C., Remer, L.A., Mattoo, S., Vermote, E.F., Kaufman, Y.J., 2007.

477        Second-generation operational algorithm: Retrieval of aerosol properties over land

478        from inversion of Moderate Resolution Imaging Spectroradiometer spectral

479        reflectance. J. Geophys. Res. Atmos. 112.

480    Liu, Y., Park, R.J., Jacob, D.J., Li, Q., Kilaru, V., Sarnat, J.A., 2004. Mapping annual mean

481        ground-level PM2. 5 concentrations using Multiangle Imaging Spectroradiometer

482        aerosol optical thickness over the contiguous United States. J. Geophys. Res. Atmos.

483        109.

484    Lu, X., Wang, J., Yan, Y., Zhou, L., Ma, W., 2021. Estimating hourly PM2. 5

485        concentrations using Himawari-8 AOD and a DBSCAN-modified deep learning model

486        over the YRDUA, China. Atmos. Pollut. Res. 12, 183–192.

487    Ma, Z., Hu, X., Sayer, A.M., Levy, R., Zhang, Q., Xue, Y., Tong, S., Bi, J., Huang, L., Liu,

488        Y., 2016. Satellite-based spatiotemporal trends in PM2. 5 concentrations: China,

489        2004–2013. Environ. Health Perspect. 124, 184–192.

490    Mahasakpan, N., Chaisongkaew, P., Inerb, M., Nim, N., Phairuang, W., Tekasakul, S.,

491        Furuuchi, M., Hata, M., Kaosol, T., Tekasakul, P., 2023. Fine and ultrafine particle-

492        and gas-polycyclic aromatic hydrocarbons affecting southern Thailand air quality

493        during transboundary haze and potential health effects. J. Environ. Sci. 124, 253–267.

494    Muñoz Sabater, J., 2019. ERA5-Land hourly data from 1981 to present. Copernicus Clim.

495        Chang. Serv. Clim. Data Store 10.

496    Organization, W.H., 2016. Ambient air pollution: A global assessment of exposure and

497        burden of disease.

498    Peters, B., Niculae, V., Martins, A.F.T., 2019. Sparse sequence-to-sequence models. arXiv

499        Prepr. arXiv1905.05702.

500    Pinichka, C., Makka, N., Sukkumnoed, D., Chariyalertsak, S., Inchai, P., Bundhamcharoen,

501        K., 2017. Burden of disease attributed to ambient air pollution in Thailand: A GIS-

502        based approach. PLoS One 12, e0189909.

503    Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A.V., Gulin, A., 2018. CatBoost:

504        unbiased boosting with categorical features. Adv. Neural Inf. Process. Syst. 31.

505    Pu, Q., Yoo, E.-H., 2021. Ground PM2. 5 prediction using imputed MAIAC AOD with

506        uncertainty quantification. Environ. Pollut. 274, 116574.

507    Punsompong, P., Pani, S.K., Wang, S.-H., Pham, T.T.B., 2021. Assessment of biomass-

508        burning types and transport over Thailand and the associated health risks. Atmos.

509        Environ. 247, 118176.

510 Rodríguez-Urrego, D., Rodríguez-Urrego, L., 2020. Air quality during the COVID-19:

511      PM2. 5 analysis in the 50 most polluted capital cities in the world. Environ. Pollut.

512      266, 115042.

513 Shin, M., Kang, Y., Park, S., Im, J., Yoo, C., Quackenbush, L.J., 2020. Estimating ground-

514      level particulate matter concentrations using satellite-based data: a review. GIScience

515      Remote Sens. 57, 174–189.

516 Sorek-Hamer, M., Strawa, A.W., Chatfield, R.B., Esswein, R., Cohen, A., Broday, D.M.,

517      2013. Improved retrieval of PM2. 5 from satellite data products using non-linear

518      methods. Environ. Pollut. 182, 417–423.

519 Stratoulias, D., Nuthammachot, N., 2020. Air quality development during the COVID-19

520      pandemic over a medium-sized urban area in Thailand. Sci. Total Environ. 746,

521      141320.

522 Theys, N., De Smedt, I., Yu, H., Danckaert, T., van Gent, J., Hörmann, C., Wagner, T.,

523      Hedelt, P., Bauer, H., Romahn, F., 2017. Sulfur dioxide retrievals from TROPOMI

524      onboard Sentinel-5 Precursor: algorithm theoretical basis. Atmos. Meas. Tech. 10,

525      119–153.

526 Tucker, W.G., 2000. An overview of PM2. 5 sources and control strategies. Fuel Process.

527      Technol. 65, 379–392.

528 Van Donkelaar, A., Martin, R. V, Brauer, M., Kahn, R., Levy, R., Verduzco, C.,

529      Villeneuve, P.J., 2010. Global estimates of ambient fine particulate matter

530      concentrations from satellite-based aerosol optical depth: development and

531      application. Environ. Health Perspect. 118, 847–855.

532 Van Geffen, J., Eskes, H.J., Boersma, K.F., Maasakkers, J.D., Veefkind, J.P., 2019.

533      TROPOMI ATBD of the total and tropospheric NO2 data products. DLR Doc.

534    Vassanadumrongdee, S., Matsuoka, S., 2005. Risk perceptions and value of a statistical life

535        for air pollution and traffic accidents: evidence from Bangkok, Thailand. J. Risk

536        Uncertain. 30, 261–287.

537    Wang, Y., Yuan, Q., Li, T., Tan, S., Zhang, L., 2021. Full-coverage spatiotemporal

538        mapping of ambient PM2. 5 and PM10 over China from Sentinel-5P and assimilated

539        datasets: Considering the precursors and chemical compositions. Sci. Total Environ.

540        793, 148535.

541    Wei, J., Huang, W., Li, Z., Xue, W., Peng, Y., Sun, L., Cribb, M., 2019. Estimating 1-km-

542        resolution PM2. 5 concentrations across China using the space-time random forest

543        approach. Remote Sens. Environ. 231, 111221.

544    Weichenthal, S.A., Godri Pollitt, K., Villeneuve, P.J., 2013. PM2. 5, oxidant defence and

545        cardiorespiratory health: a review. Environ. Heal. 12, 1–8.

546    Wong, P.-Y., Lee, H.-Y., Chen, Y.-C., Zeng, Y.-T., Chern, Y.-R., Chen, N.-T., Lung, S.-

547        C.C., Su, H.-J., Wu, C.-D., 2021. Using a land use regression model with machine

548        learning to estimate ground level PM2. 5. Environ. Pollut. 277, 116846.

549    World Air Quality Report 2020 [WWW Document], n.d. URL

550        https://www.greenpeace.org/static/planet4-romania-stateless/2021/03/d8050eab-2020-

551        world_air_quality_report.pdf

552    World Economic Outlook (October 2021) [WWW Document], n.d. URL

553        https://www.imf.org/external/datamapper/datasets/WEO (accessed 2.7.22).

554    Xie, Y., Wang, Y., Zhang, K., Dong, W., Lv, B., Bai, Y., 2015. Daily estimation of ground-

555        level PM2. 5 concentrations over Beijing using 3 km resolution MODIS AOD.

556        Environ. Sci. Technol. 49, 12280–12288.

557    Xue, T., Zheng, Y., Tong, D., Zheng, B., Li, X., Zhu, T., Zhang, Q., 2019. Spatiotemporal

558         continuous estimates of PM2. 5 concentrations in China, 2000–2016: A machine

559         learning method with inputs from satellites, chemical transport model, and ground

560         observations. Environ. Int. 123, 345–357.

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

562         ground-level PM2. 5 concentration in China using geographically weighted regression

563         based on 3 km resolution MODIS AOD. Remote Sens. 8, 184.

564   Zhang, X., Han, Liangxiu, Han, Lianghao, Zhu, L., 2020. How well do deep learning-based

565         methods for land cover classification and object detection perform on high resolution

566         remote sensing imagery? Remote Sens. 12, 417.

567   Zhu, L., Huang, L., Fan, L., Huang, J., Huang, F., Chen, J., Zhang, Z., Wang, Y., 2020.

568         Landslide susceptibility prediction modeling based on remote sensing and a novel

569         deep learning algorithm of a cascade-parallel recurrent neural network. Sensors 20,

570         1576.

571   Zou, B., Chen, J., Zhai, L., Fang, X., Zheng, Z., 2017. Satellite based mapping of ground

572         PM2. 5 concentration using generalized additive modeling. Remote Sens. 9, 1.

573