# A Satellite-based Model for Estimating PM$_{2.5}$ Concentration in a Sparsely Populated Environment Using Soft Computing Techniques

*Bijan Yeganeh[1, 2], Michael G. Hewson[3], Samuel Clifford[2, 4], Luke D. Knibbs[5], and Lidia Morawska[1, \*]*

[1] International Laboratory for Air Quality and Health, Queensland University of Technology, Brisbane, Queensland 4001, Australia

[2] Centre for Air Quality and Health Research and Evaluation, Glebe, New South Wales 2037, Australia

[3] School of Education and the Arts, Central Queensland University, North Rockhampton, Queensland 4702, Australia

[4] Mathematical Sciences School, Queensland University of Technology, Brisbane, Queensland 4000, Australia

[5] School of Public Health, The University of Queensland, Herston, Queensland 4006, Australia

**\* Corresponding author:** Professor Lidia Morawska, PhD

International Laboratory for Air Quality and Health

Queensland University of Technology, Brisbane, Queensland 4001, Australia

Phone: +61 7 3138 2616

Fax: +61 7 3138 9079

Email: l.morawska@qut.edu.au

**Highlights**

- We used comprehensive satellite-based, meteorological and geographical data to develop a satellite-based model for estimating the $PM_{2.5}$ concentration.

- Representative animations are created to visualize the spatiotemporal variation of the predictors.

- We applied the adaptive neuro-fuzzy inference system (ANFIS) for the first time as a core model to estimate the spatiotemporal variation of $PM_{2.5}$ concentration.

- We compared ANFIS with support vectors machine and back-propagation artificial neural network.

- Adaptive model identification technique has been used to identify the optimal predictive model.

44    **Abstract**

45    We applied three soft computing methods including adaptive neuro-fuzzy inference system

46    (ANFIS), support vectors machine (SVM) and back-propagation artificial neural network

47    (BPANN) algorithms for estimating the ground-level $PM_{2.5}$ concentration. These models were

48    trained by comprehensive satellite-based, meteorological, and geographical data. A 10-fold

49    cross-validation (CV) technique was used to identify the optimal predictive model. Results

50    showed that ANFIS was the best-performing model for predicting the variations in $PM_{2.5}$

51    concentration. Our findings demonstrated that the CV-$R^2$ of the ANFIS (0.81) is greater than that

52    of the SVM (0.67) and BPANN (0.54) model. The results suggested that soft computing methods

53    like ANFIS, in combination with spatiotemporal data from satellites, meteorological data and

54    geographical information improve the estimate of $PM_{2.5}$ concentration in sparsely populated

55    areas.

56

57

58    **Keywords**: $PM_{2.5}$; Aerosol optical depth; ANFIS; SVM; BPANN; Australia.

59

60

61

62

63 **Data availability**

64

65 The type and source of the data set considered in this study.

| Name of the data set | Data source (Developer) (All websites accessed on Jan 2016) | Data format | Software required | Data availability |
|---|---|---|---|---|
| **OMI Near-UV AOD** | Aura OMI AOD product via NASA Giovanni interface http://giovanni.sci.gsfc.nasa.gov/giovanni/?instance_id=omil2g | HDF / NetCDF files | ArcGIS | Freely available |
| **Major road** | PSMA Australia Transport and Topography product https://www.psma.com.au/products/transport-topography | ESRI shape files | " " | Price depends on the area of interest |
| **Minor road** | " " | " " | " " | " " |
| **Industrial point source PM$_{2.5}$ emissions** | Australia National Pollutant Inventory http://www.npi.gov.au/reporting/industry-reporting-materials | xml files | Microsoft Excel / R | Freely available |
| **Australia population density** | Australian Bureau of Statistics http://www.abs.gov.au/ausstats/abs@.nsf/mf/1270.0.55.007 | PNG ESRI Grid GeoTIFF | ArcGIS | " " |
| **Australia land use classification** | Australian Bureau of Statistics http://www.abs.gov.au/websitedbs/censushome.nsf/home/meshblockcounts | Excel spreadsheets / CSV files | Microsoft Excel / R / ArcGIS | " " |
| **Elevation** | U.S. Geological Survey https://www.usgs.gov/products/maps/topo-maps | PNG GeoTIFF | ArcGIS | " " |
| **Normalized difference vegetation index** | Terrestrial Ecosystem Research Network http://www.auscover.org.au/node/9 | NetCDF files | " " | " " |
| **Temperature** **Rainfall** **Humidity** **Solar exposure** | Australian Bureau of Meteorology http://www.bom.gov.au/climate/maps/#tabs=Maps | ESRI Grid GIF | " " | " " |

66

67 **Software availability**

68 The following software has been used in this study for statistical analysis, spatial data processing

69 and map creation:

70 - R v.3.2.3 (R Foundation for Statistical Computing, Vienna, Austria)

71 - MATLAB R2014b (MathWorks Inc., Natick, USA)

72 - ArcGIS version 10.2 (ESRI Inc., Redlands, USA)

73 **Note:** No specific software component has been developed for this study.

74

75

76

77

78

79

80

81

82

## 1. Introduction

Exposure to fine particulate matter (PM$_{2.5}$, particles with aerodynamic diameter less than 2.5 μm) is a leading environmental risk factor associated with respiratory and cardiovascular morbidity and mortality (Franklin et al., 2007) and it is the twelfth-ranked contributor to the global burden of diseases (Forouzanfar et al., 2015).

Urbanisation increases the risk of being exposed to PM$_{2.5}$ (Han et al., 2015), and Australia, as one of the most urbanised countries in the world , is faced with adverse health effects of PM$_{2.5}$. To date, very little attention has been paid to the health effect of exposure to PM$_{2.5}$ in Australia. Some studies consistently suggest that PM$_{2.5}$ is associated with respiratory diseases and has significant effects on mortality (Barnett et al., 2005; Simpson et al., 2005), while conflicting results have been reported on cardiovascular health effects (Hinwood et al., 2006). These inconsistent results could be due to difficulties in assessing the Australian population exposure to PM$_{2.5}$.

Ground level aerosol measurement has been historically provided by ground monitoring networks, but there are high establishing and maintaining expenses associated with these measurements (Wu et al., 2012). The sparse ground PM$_{2.5}$ measurement network in Australia makes it difficult to evaluate the spatiotemporal variability of PM$_{2.5}$ and has significantly restrained the epidemiological studies on PM$_{2.5}$ health effects. Australia is the sixth largest country in the world by area while its population is quite small compared to the land size (Australian Government, 2015). Australia is one of the 10 least dense populated countries in the world (United Nations, 2015). The majority of the Australian population is living in the east and west coasts (Lunn et al., 2002). The population within these areas is concentrated in urban centres, particularly the capital cities (Australian Bureau of Statistics, 2012; Lunn et al., 2002).

106 Therefore, limited monitoring stations were established only in populated areas due to population

107 distribution in Australia. Had such monitoring networks existed, there would have been no

108 guarantee of an effective measurement of the spatiotemporal variation of $PM_{2.5}$, since it is

109 changing on scales much smaller than monitoring networks density.

110 Estimates of air pollution exposure have been traditionally provided by assigning

111 measurements derived from one (Chen et al., 2006) or several air pollution monitors (Barnett et

112 al., 2005; Brook et al., 2010; Chan et al., 2006), allocating exposure using the nearest monitoring

113 station (Lee et al., 2014) or using different proxies to estimate a local population's exposure

114 (Hoffmann et al., 2007; Salam et al., 2008; Samet, 2007). There is potential for over-smoothing

115 the exposure estimation and the results are likely to be biased with all these approaches (Jerrett et

116 al., 2005a).

117 Satellite imagery is another important tool rapidly gaining interest in air pollution monitoring

118 as it provides sequential observations over a broad area. Satellite sensors can be coupled with

119 ground-based sensors at different spatiotemporal scales to reduce the limitations of surface

120 monitoring station (Reis et al., 2015). Aerosol Optical Depth (AOD) is the most common

121 parameter derived from satellite observations and applied to estimate $PM_{2.5}$. AOD describes the

122 level of which aerosols attenuate the electromagnetic radiation at a given wavelength by

123 absorption or scattering in an atmospheric column (Chudnovsky et al., 2012; Kaufman et al.,

124 2002; NASA, 2013). The availability of satellite-derived AOD has helped to overcome the

125 problems associated with sparse monitoring networks by providing observations where

126 previously there were none (Hoff and Christopher, 2009; Reis et al., 2015).

127 A variety of methods have been used to investigate the quantitative relationship between

128 satellite-derived AOD and ground-level $PM_{2.5}$ measurements. These studies mainly fall into two

129    major classes: numerical-based methods and empirical observation-based methods (Lin et al.,

130    2014).

131    Numerical-based models, including dispersion and chemical transport models, are still under

132    development due to the uncertainties regarding the definition of source inventories, and chemical

133    and dynamical processes of aerosols in atmosphere (Gupta and Christopher, 2009b; Kondragunta

134    et al., 2008). Empirical observation-based methods rely on the relationship between air quality

135    measurements and different observations (Maciejewska et al., 2015). Several techniques have

136    been used to describe this relationship including simple regression (Chu et al., 2003), multiple

137    regression (Dirgawati et al., 2015; Gupta and Christopher, 2009b; Li et al., 2011), geostatistical

138    methods (Jerrett et al., 2005b; Kunzli et al., 2005), generalized additive models (GAM) (Strawa

139    et al., 2013), land use regression (Henderson et al., 2007; Kloog et al., 2011; Knibbs et al., 2014;

140    Liu et al., 2009), and hybrid approaches (Beckerman et al., 2013b; Lindstrom et al., 2011). Soft

141    computing refers to computational techniques which are able to achieve optimal solutions for

142    analysing complicated phenomena at reasonable costs (Carnevale et al., 2016; Kruse et al., 2013;

143    Ovaska, 2004). In recent years, soft computing techniques such as support vector machine

144    (SVM) (Moazami et al., 2016; Reid et al., 2015; Yeganeh et al., 2012), Bayesian models (Corani

145    and Scanagatta, 2016; McBride et al., 2007), $k$-nearest neighbours (kNN) (Reid et al., 2015), and

146    artificial neural network (ANN) (Al-Alawi et al., 2008; Gupta and Christopher, 2009a; Ordieres

147    et al., 2005; Wu et al., 2012) have been gaining popularity in air quality modelling because of

148    their high flexibility and well documented prediction abilities. However, other soft  comuting

149    methodes like adaptive neuro-fuzzy inference system (ANFIS), which is accepted as an efficient

150    and robust method for multivariate analysis, have not been used for modelling the spatiotemporal

151    variations of $PM_{2.5}$ concentrations.

152　　Although most of the aforementioned methods can be applied to determine the relationship

153　　among AOD and PM$_{2.5}$, imposing a specific method could make it difficult to select the best

154　　predictive model. Hence, adaptive model identification approach is used to choose the most

155　　efficient model by using cross-validation technique rather than fitting a specific model to the

156　　dataset (Reid et al., 2015; Syed, 2011).

157　　Few studies have investigated the relationship between PM$_{2.5}$ and satellite-based AOD in

158　　Australia (Gupta et al., 2007; Gupta et al., 2006; Meyer et al., 2008). While other studies have

159　　recommended the meteorological and geographical factors incorporation to the AOD–PM$_{2.5}$

160　　relationship to improve models' performance (Chudnovsky et al., 2014; Liu et al., 2009), there is

161　　a clear need to conduct an Australian study to develop a satellite-based model investigating

162　　significant geographical and meteorological factors including humidity, planetary boundary

163　　layer, and wind speed and direction.

164　　In this study, we aimed to improve the estimate of PM$_{2.5}$ concentration by using remotely-

165　　sensed AOD in conjunction with comprehensive meteorological and geographical data.　Three

166　　different soft computing algorithms were applied to estimate the monthly average exposure to

167　　PM$_{2.5}$ in the South-east Queensland (SEQ) region of Australia, from 2006 to 2011. In turn, an

168　　adaptive model identification approach was used to choose the optimal model from ANFIS and

169　　other soft computing methods: SVM and BPANN, by using 10-fold cross-validation (Pandey et

170　　al., 2013; Syed, 2011). We ultimately used the model with the best predictive ability to estimate

171　　spatiotemporal variation of PM$_{2.5}$ in this sparsely populated area with dense vegetation cover.

172　　**2. Materials and Methods**

173　　**2.1. Study location and ground-level PM$_{2.5}$ measurements**

174    SEQ is a region in the state of Queensland, Australia, which covers 22,420 km$^2$ and is home to

175    3.05 million people out of the state's population of 4.58 million based on the 2011 Australian

176    census (Australian Bureau of Statistics., 2012). The study area consists of Brisbane, the state's

177    capital city, as well as other urban and rural centres including Ipswich, Logan City, Gold Coast,

178    Sunshine Coast, and the Lockyer Valley. Motor vehicle emissions and industrial boilers are

179    identified as major sources of PM$_{2.5}$ in SEQ (Queensland Government., 2014). The Queensland

180    government and other agencies are responsible for regulatory aerosol monitoring in SEQ. We

181    obtained quality-assured 24 h ground-level PM$_{2.5}$ measurements from January 2006 to December

182    2011. During the study period, PM$_{2.5}$ was measured at 8 monitoring sites across SEQ

183    (supplement, page S3). We used monthly averages of the daily measured PM$_{2.5}$, and the

184    inclusion criteria for a given month was that less than 5% of the daily measurements were

185    missing.

186

187    **2.2. Land use data**

188    We obtained data on anthropogenic and natural land use variables as spatial predictors that

189    were possible predictors of measured PM$_{2.5}$ concentrations. The selected land use variables,

190    summarised in Table 1, were examined to discover which ones improved the prediction of PM$_{2.5}$

191    (Knibbs et al., 2014). They included proxies for emissions from traffic, point sources and

192    changing land cover conditions.

193    The impacts of vegetation cover and its phenological state on the relationship between the

194    PM$_{2.5}$ and satellite AOD were also examined in the present study. Normalized difference

195    vegetation index (NDVI) is used to provide a measure of greenness and vegetation cover. NDVI

196  was found to be an effective predictor for pollutant concentrations in previous studies

197  (Chudnovsky et al., 2014; Dirgawati et al., 2015; Su et al., 2009). The monthly mean NDVI data

198  were derived from an Advanced Very High Resolution Radiometer (AVHRR) sensor carried on

199  the National Oceanic and Atmospheric Administration (NOAA) satellite and processed by the

200  Australian Bureau of Meteorology (BoM) at a spatial resolution of 1 km (Bureau of

201  Meteorology, 2015).

202  **2.3. Satellite data**

203  Daily global gridded observations of AOD at a resolution of 0.25 degrees latitude and

204  longitude are derived from the Ozone Monitoring Instrument (OMI) aboard the Aura satellite

205  (Levelt et al., 2006). Aura crosses the equator in a sun-synchronous polar orbit for the daylight

206  ascending orbit (Torres et al., 2007), and it passes over SEQ at approximately 14:00 local time.

207  We download the monthly average OMI AOD level 2 Near-UV AOD and single Scattering

208  Albedo product (OMAERUVG.003 at 342.5 nm) from NASA Giovanni interface for each month

209  from 2006–2011.

210

211

212

213

214

215

216

217

218

219    Table 1. Independent variables included as potential predictors of PM$_{2.5}$

| Variables (units) | Spatial resolution | Point or buffer | Data source |
|---|---|---|---|
| **OMI Near-UV AOD** | 0.25 degrees Lat/Lon | Point | Aura OMI AOD product via NASA Giovanni interface |
| **Distance to coast (m)** | - | Point | ArcGIS geoprocessing tools |
| **Distance to port (m)** | - | Point | " " |
| **Distance to airport (m)** | - | Point | " " |
| **Distance to nearest major road** | - | Point | " " |
| **Distance to nearest minor road** | - | Point | " " |
| **Airport (present/not present)** | - | Buffer | " " |
| **Major road (m)** | - | Buffer | PSMA Australia Transport and Topography product |
| **Minor road (m)** | - | Buffer | " " |
| **Industrial point source PM$_{2.5}$ emissions (kg/yr)** | - | Buffer | Australia National Pollutant Inventory |
| **Time (Julian month)** | - | Point | ArcGIS geoprocessing tools |
| **Population density (person/km$^2$)** | $1 \times 1$ km$^2$ | Point | Australian Bureau of Statistics |
| **Land use by type (% area)[b]** | Mesh block[c] | Buffer | Australian Bureau of Statistics |
| **Elevation (m)** | 30 m | Point | U.S. Geological Survey |
| **Normalized difference vegetation index** | $1 \times 1$ km$^2$ | Point | Terrestrial Ecosystem Research Network, Australian Bureau of Meteorology, AusCover project and NASA NOAA satellite |
| **Mean daily maximum temperature (°C)** | $5 \times 5$ km$^2$ | Point | Australian Bureau of Meteorology |
| **Mean daily minimum temperature (°C)** | $5 \times 5$ km$^2$ | Point | " " |
| **Rainfall (mm)** | $5 \times 5$ km$^2$ | Point | " " |
| **Humidity (hPa)** | $5 \times 5$ km$^2$ | Point | " " |
| **Solar exposure (MJ/m$^2$)** | $6 \times 6$ km$^2$ | Point | " " |
| **Planetary boundary layer height (m)** | $3 \times 3$ km$^2$ | Point | Derived from Weather Research and Forecasting model |
| **U-component of wind speed (m/s)** | $3 \times 3$ km$^2$ | Point | " " |
| **V-component of wind speed (m/s)** | $3 \times 3$ km$^2$ | Point | " " |
| **Wind speed (m/s)** | $3 \times 3$ km$^2$ | Point | " " |
| **Wind direction (Degrees)** | $3 \times 3$ km$^2$ | Point | " " |

220
221    [a] 22 Circular buffers were generated with radii of 50 m, 100 m, 200 m, 300 m, 400 m, 500 m, 600 m, 700 m, 800 m, 900 m, 1000 m, 1200 m, 1500 m, 1800 m, 2000 m, 2500 m, 3000 m, 3500 m, 4000 m, 5000 m, 7500 m, and 10,000 m (Novotny et al., 2011).

222
223    [b] Four different land use classes were investigated including industrial, commercial, residential, and open space (which contains the agricultural land, parks, and water bodies (Rose et al., 2010)).

224
225    [c] Mesh Block is the smallest geographic unit defined by the Australian Statistical Geography Standard for which the Census data is available (Australian Bureau of Statistics, 2011), and can be variable in size.

226

227

## 2.4. Meteorological Data

228

229 We obtained surface meteorological parameters including mean maximum and minimum

230 temperature, rainfall, and humidity from high-quality spatial climate data-sets developed by

231 BoM which provides gridded climatological maps for each month of the year (Jones et al.,

232 2009). In addition, monthly solar exposure maps are also obtained from BoM during the study

233 period.

234 Planetary boundary layer height (PBLH), wind direction (WD) and wind speed (WS) can play

235 a critical role in the transport and dilution of $PM_{2.5}$ (Harrison et al., 1997); hence, special

236 attention was paid to these parameters in this study. The Weather Research and Forecasting

237 model (WRF) was used to calculate these parameters as at 2:00 pm local time at a spatial

238 resolution of 3 km to match the over-pass time of the Aura satellite. Details on the WRF

239 configuration are provided in the supplement (page S3-S7).

## 2.5. Modelling approach

240

241 Following similar studies (Knibbs et al., 2014; Novotny et al., 2011), 22 circular buffers were

242 created around each monitoring site to obtain local and more remote sources of $PM_{2.5}$ (Table 1).

243 Certain variables were calculated within a buffer (e.g., land use type, road length) while others

244 were extracted at each monitoring site (e.g., wind speed, humidity, temperature). In total, 194

245 independent variables were obtained including 18 point variables and 176 buffer variables (8

246 variables calculated at 22 buffers each).

## 2.6. Statistical analysis

247

13

248     In this study, there were 194 predictor variables to choose from, hence choosing the optimum

249    subset was a complicated process and needed to be carefully conducted. In many soft computing

250    and data mining tasks, there can be some irrelevant variables which may affect the derived

251    statistical relationship between the dependent variable and the other relevant predictors. A

252    common solution to overcome this problem is to use a variable selection process which can help

253    to select a subset of the most relevant and representative predictors from input predictors. The

254    Least Absolute Shrinkage and Selection Operator (Lasso) is a well-known method which is

255    widely used to suppress or shrink variables to select the most relevant predictor variable set.

256    Lasso-type variable selection method was used in this study since it was successfully adopted in

257    many applications (Hu et al., 2015; Li and Shao, 2015; Tibshirani, 2011).

258    Following Beelen et al. (2013), we only included potential predictor variables with less than

259    10% null values and centred and standardised some independent variables to improve model

260    convergence and make the parameter estimates more interpretable. Subsequently, all remaining

261    predictor variables were evaluated, and variables with p-value greater than 0.10 or variance

262    inflation factor (VIF) greater than 6 were removed in order to avoid multicollinearity (see Table

263    S2). As suggested by other studies (Beelen et al., 2013; Henderson et al., 2007; Novotny et al.,

264    2011; Vienneau et al., 2013), if two buffer sizes of a particular variable were found to be

265    collinear, donut ring buffers (so called concentric adjacent rings) were replaced with original

266    circular buffers and the analysis was redone. Ring buffers (i.e., annulus) were calculated by

267    differencing the circular buffers.

268    In this study, we employed the soft computing techniques ANFIS, BPANN and SVM. The soft

269    computing techniques employed here are explained in the supplement, page S8-S17. The input

270    variables were composed of different types of data, including land use, meteorological, and
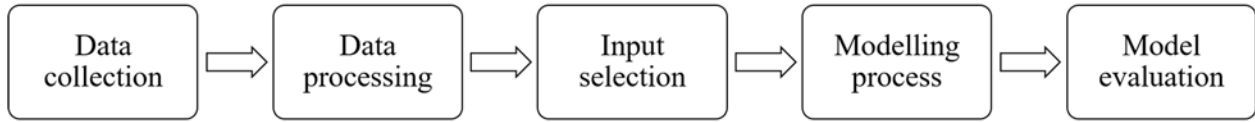
271     satellite data. We matched the selected variables with the $PM_{2.5}$ measurements at 8 monitoring

272     sites during the study period (72 months) which resulted in more than 550 observation sets in

273     total. These observations were divided into training, validation, and test subsets. The majority of

274     the observations (70%) were used for training the models. In order to avoid over-training, 15%

275     of the observations were used for validation and checking the model's generalisation (Wu et al.,

276     2012). Finally, the remaining 15% of the observations were employed as the test subset to

277     estimate the $PM_{2.5}$ concentration by the models.

278     In this study, 10-fold cross validation (CV) method was applied to evaluate the performance of

279     the BPANN, ANFIS, and SVM models and identify the optimal model for estimating the $PM_{2.5}$

280     concentration. This method has the ability to examine the model's predictive ability (Beckerman

281     et al., 2013a). This examination was accomplished by randomly splitting the data into 10 equal-

282     sized folds. Subsequently, one of the folds was used to test the model and the remaining 9 folds

283     were used to train the model (Kim, 2009; Refaeilzadeh et al., 2009). This process was repeated

284     10 times for each candidate model while all folds were used as the test subset and the 10 results

285     were averaged to obtain the overall $CV\text{-}R^2$ and CV-RMSE. The best predictive model was

286     selected from those with the smallest CV-RMSE and highest $CV\text{-}R^2$ (Dirgawati et al., 2015).

287     Bland-Altman plot was also used to examine the agreement between the observations and

288     predictions. In this plot, X axis shows the average of the model predictions and observations, and

289     Y axis represents the difference between these values. Bland-Altman plot also provides statistical

290     limits by calculating the average and mean and the standard deviation (*sd*) of the differences

291     between observations and predictions (Giavarina, 2015). These limits were used to evaluate the

292     agreement between observations and model predictions. For more explicit information on Bland-

293    Altman plot see Giavarina (2015). Figure 1 illustrates the overall research process used in this

294    study.

295



296    Figure 1. General research process for estimating PM$_{2.5}$ concentration.

297    We used R v.3.2.3 (R Foundation for Statistical Computing, Vienna, Austria) and MATLAB

298    R2014b (MathWorks Inc., Natick, USA) for all statistical and soft computing analyses and

299    ArcGIS version 10.2 (ESRI Inc., Redlands, USA) for spatial data processing and map creation.

300    **3. Results**

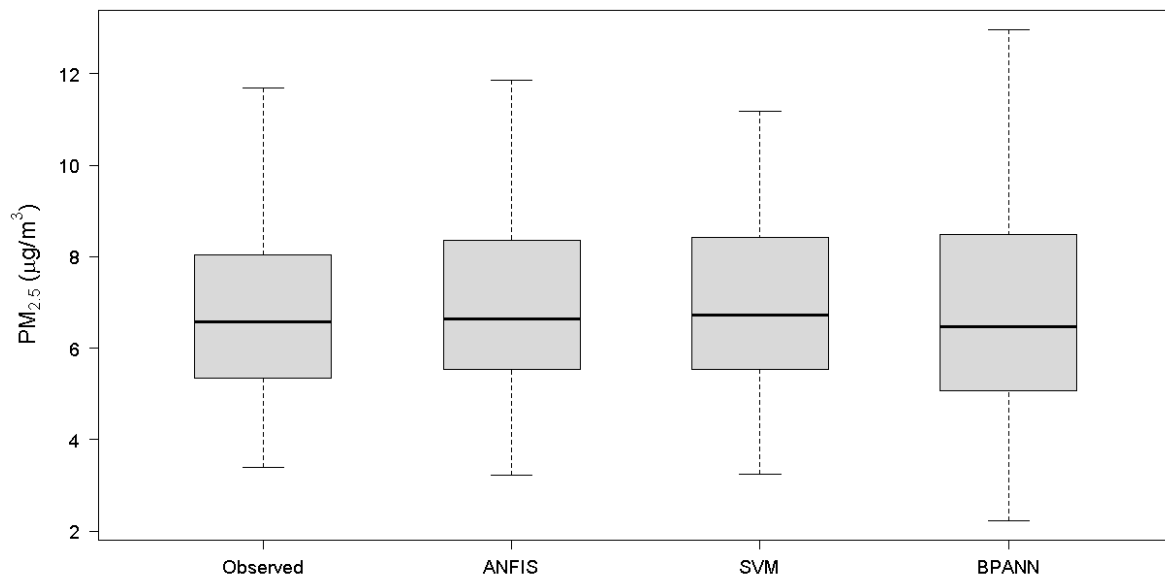301    **3.1. Modeling results and evaluation**

302    In this study, a wide range of ground-based PM$_{2.5}$ measurements, land use, meteorological, and

303    remotely-sensed AOD data were employed to estimate the PM$_{2.5}$ concentration using soft

304    computing techniques. In following section, the agreement between predicted and observed

305    PM$_{2.5}$ concentration is evaluated. 10-fold cross validation is also used to compare the potential of

306    different algorithms for estimating PM$_{2.5}$ concentration.

307    The variable selection results showed that 16 variables were the most effective predictors of

308    PM$_{2.5}$ concentration. The variables most correlated with the outcome were firstly humidity,

309    followed by maximum temperature, AOD and then the length of the major roads. The results of

310    the variable selection process are provided in the supplement (Table S2).

311    Using the testing dataset for each of the developed models, PM$_{2.5}$ concentrations were then

312    predicted. A summary of the observed and predicted PM$_{2.5}$ concentrations is presented in Fig. 2.

16

313    The mean observed PM$_{2.5}$ concentration for the testing dataset is 6.77 µg/m$^3$. All three models

314    approached this value within a numerical range of -0.02 to +0.38 µg/m$^3$. The non-parametric

315    Wilcoxon test was performed to check if there was any significant difference between the

316    predicted and observed mean PM$_{2.5}$ concentrations of each model. The test on all models yielded

317    p-values greater than 0.01, showing an insignificant difference between the predicted and

318    observed PM$_{2.5}$ concentration at 1% significant level. A comparison of the predicted values

319    demonstrated that the ANFIS model predicted values were slightly closer to the full range of the

320    observed monitoring data than the SVM and BPANN models. In general, Fig. 2 shows that all

321    models reliably calculated the average and range of PM$_{2.5}$ concentration; therefore, predicted-

322    observed plots are used to evaluate the predictive abilities of the models (Fig. 3).
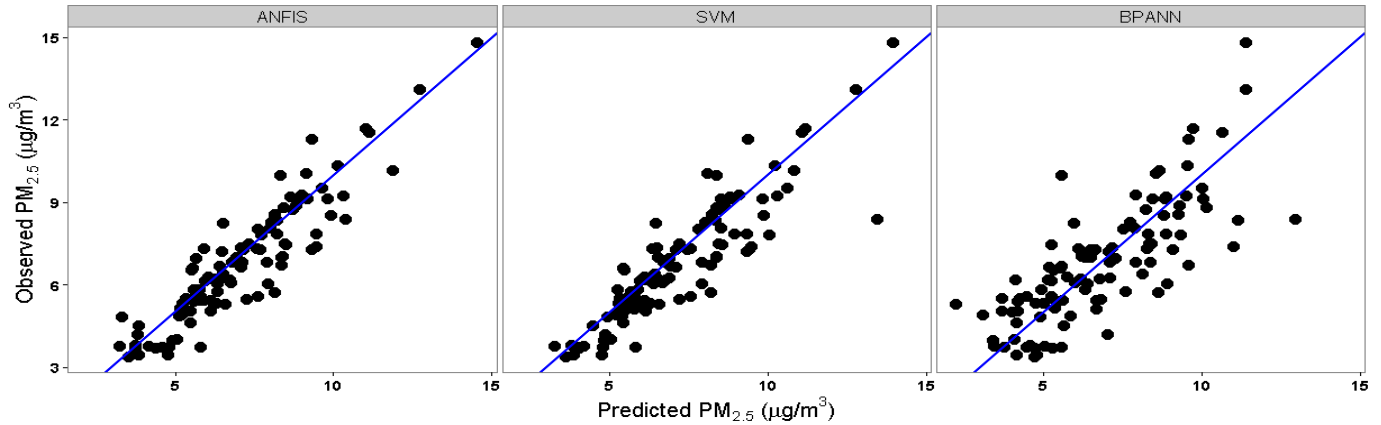
323



324

325    Figure 2. Summary of PM$_{2.5}$ model predictions on the testing dataset.

326    We compared the observed PM$_{2.5}$ concentrations to the predicted values of the ANFIS, SVM,

327    and BPANN models. The ANFIS's predicted-observed plot indicates that the values are more

328    equally scattered across the line of agreement at the low and high PM$_{2.5}$ concentrations whereas

329     the SVM model under-predicts and over-predicts these values, respectively. In addition, the

330     predicted-observed plot shows relatively weak correlation between the BPANN's predictions and

331     actual observations.



332

333     Figure 3. Scatter plots of observed vs. predicted PM2.5 for the optimal model fitting on the testing
334     dataset using ANFIS, SVM, and BPANN, respectively. Blue line indicates the line of agreement
335     (y = x).

336         Table 2 compares the $R^2$ and RMSE for model fitting and cross validation. For the model fit

337     the $R^2$ values are 0.61, 0.73, and 0.84 for the BPANN, SVM and ANFIS models, respectively.

338     The RMSE values are 1.57 $\mu g/m^3$, 1.36 $\mu g/m^3$, and 0.94 $\mu g/m^3$ for the BPANN, SVM, and

339     ANFIS models, respectively. Comparing the model fittings and cross validation, CV-$R^2$

340     decreases by just 0.03 and CV-RMSE increases by 0.85 $\mu g/m^3$ for the ANFIS model indicating

341     negligible model overfit. The CV-$R^2$ of the SVM and BPANN decreased by 0.06 and 0.07,

342     respectively indicating both models overfit more than ANFIS.
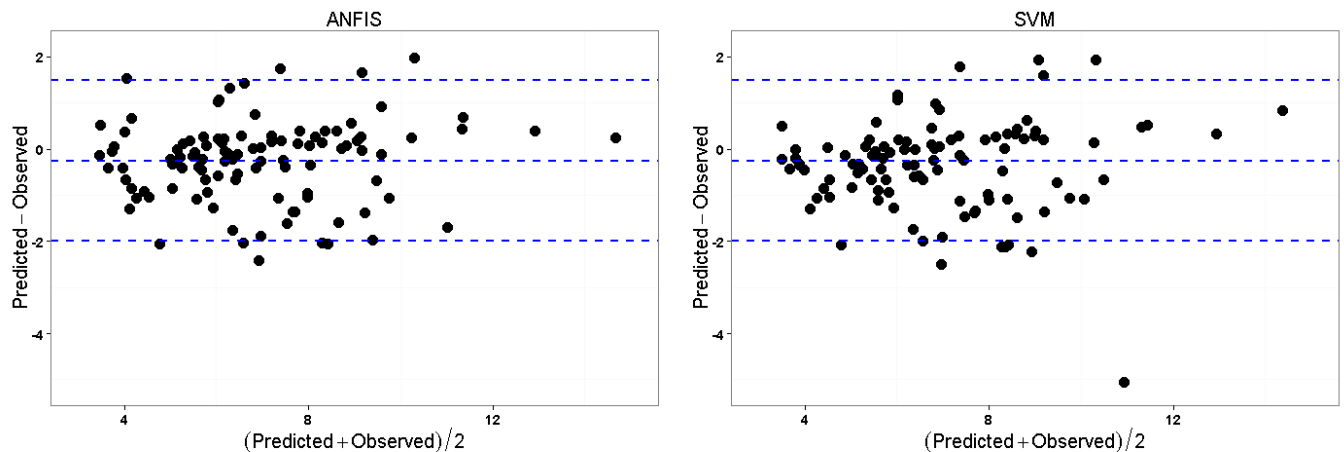
343

344

345

346

347    Table 2. R-squared and RMSE for model fittings vs. cross validation

|  | $R^2$ | RMSE ($\mu g/m^3$) | CV-$R^2$ | CV-RMSE ($\mu g/m^3$) |
|---|---|---|---|---|
| **ANFIS** | 0.84 | 0.94 | 0.81 | 1.79 |
| **SVM** | 0.73 | 1.36 | 0.67 | 2.02 |
| **BPANN** | 0.61 | 1.57 | 0.54 | 2.11 |

348

349    Our findings demonstrated that the CV-$R^2$ of the ANFIS (0.81) was higher than that of the

350    SVM (0.67) and BPANN (0.54) model. Also, the CV-RMSE of the ANFIS model (1.79 $\mu g/m^3$)

351    was lower than that of the SVM (2.02 $\mu g/m^3$) and BPANN (2.11 $\mu g/m^3$) model. Compared to

352    SVM and BPANN models, the ANFIS model had higher accuracy without causing more overfit.

353    Bland-Altman analysis was used to evaluate the agreement between the observation and

354    predictions of ANFIS and SVM since both models showed promising performance in the testing

355    stage (Figure 4). The Bland-Altman plots demonstrated low bias in both models; however, the

356    ANFIS model had slightly tighter agreement than the SVM with fewer large residuals.
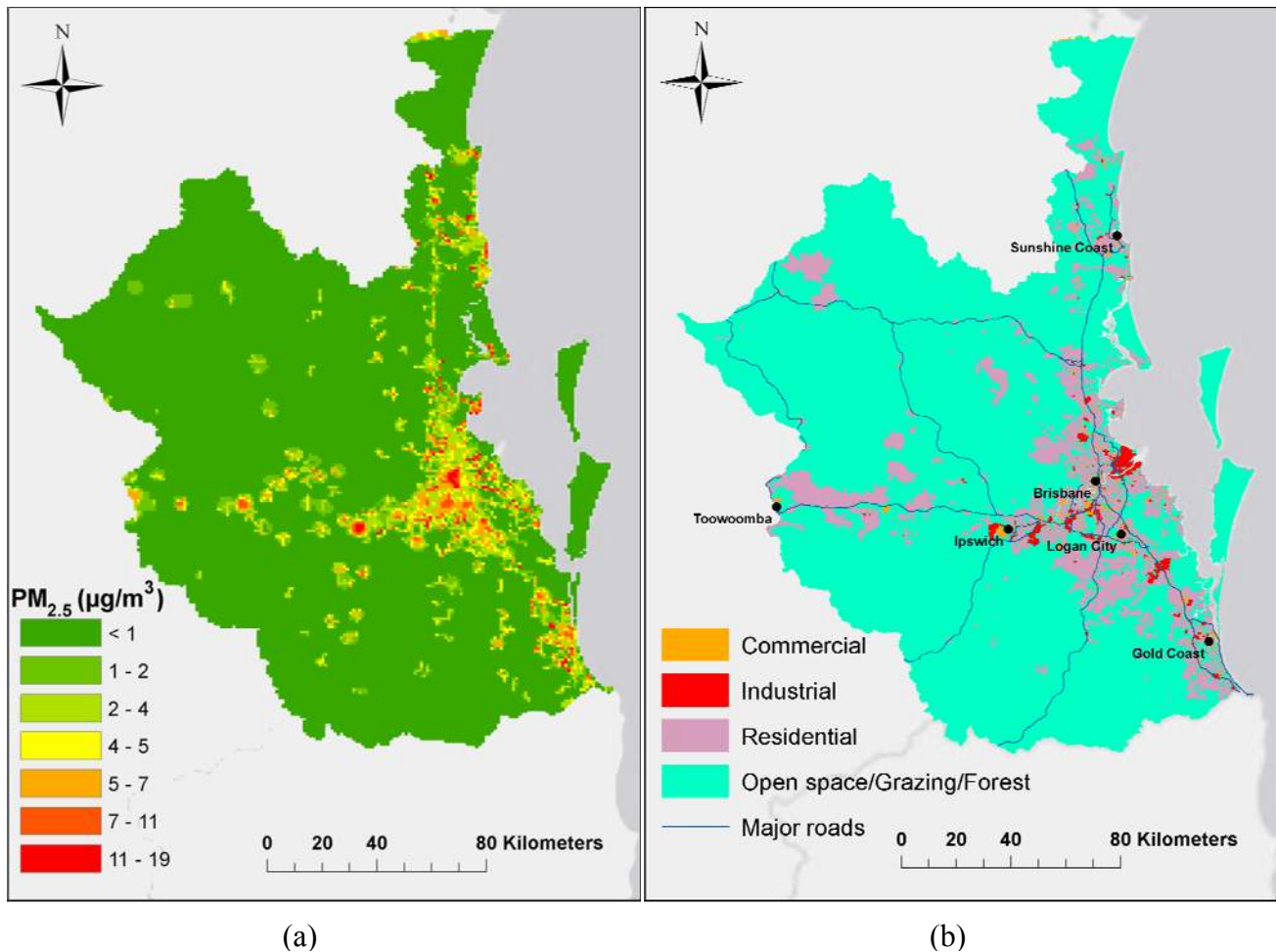


357

358    Figure 4. Bland-Altman plots of predicted and observed PM$_{2.5}$ concentrations ($\mu g/m^3$)

359

**3.2. Application of the model**

361    The predicted values of PM$_{2.5}$ concentration in September 2011, selected from the 6-year study

362    period, using ANFIS model for the study area 1 km grid is presented in Figure 5a.

363



(a)                                                                        (b)

364    Figure 5. a) Monthly average PM$_{2.5}$ concentration in September 2011 predicted by ANFIS model

365    b) land use map of SEQ

366    Figure 5b illustrates land use map of SEQ. Concentrations ranged from less than 2 to 19 μg/m$^3$.

367    Areas with higher concentrations (7 to 19 μg/m$^3$) corresponded to cities and major towns. Higher

368    concentrations were predicted in locations with extensive adjacent industrial areas and major

369    roads. This pattern was observed in all 6 cities of the study area. The highest levels were

370    predicted for the three largest cities: Brisbane (with 1.977 million people), Gold Coast (494,500)

371    and Logan (287,474).

372    **4. Discussion**

373    We employed soft computing techniques to improve concentration estimates for $PM_{2.5}$ using

374    satellite, meteorological and land use predictor variables in South-east Queensland, Australia.

375    The ANFIS model utilized in this work was the first attempt to apply it for spatiotemporal

376    modelling of $PM_{2.5}$. Using cross validation technique, the ANFIS model was found to have the

377    best performance compared to SVM and BPANN models, and better agreement with the

378    observed data. The results provide estimates of monthly $PM_{2.5}$ concentrations for SEQ from 2006

379    to 2011.

380    ANFIS is a hybrid system that combines the strengths of fuzzy logic and artificial neural

381    network (Jang, 1993; Lin and Lee, 1991), which provides a robust and accurate method for

382    predicting $PM_{2.5}$ concentration over the range of observations used in this study.

383    In this research, WRF used to calculate PBLH, and WS. Both parameters were highly

384    associated with $PM_{2.5}$ concentration across ANFIS runs. In addition, daily maximum temperature

385    had higher importance compared to the daily minimum temperature. This might be because daily

386    maximum temperature is temporally coincident with the Aura satellite overpass time for the

387    study area.

388    Prior studies mostly considered land use parameters mainly focused on roadways and

389    population-related data. We also evaluated industrial point source emissions, port and airport

390    locations as potential land use predictors. Land use parameters used in our study, may not be

391     important predictors for short term events (e.g. bush fire episode) but our findings revealed that

392     they are of the strong predictors for $PM_{2.5}$ estimation in a long term period.

393     Data sets with different spatial resolutions have been used in our study. The resolution of NDVI

394     and WRF outputs for example are finer than the OMI sensor data. Individually, OMI AOD data

395     are not spatially fine enough for estimating the $PM_{2.5}$ exposure in epidemiological studies.

396     However, method of combining different buffer sizes of land use parameters with meteorological

397     and satellite-based data enabled the model to integrate fine and more spatially coarse data sets to

398     estimate $PM_{2.5}$ concentration and provide more informative results for epidemiological studies.

399     Our results also corroborated with Reis et al. (2015) who suggested that the incorporation of 'big

400     data' derived from different sources provides new opportunities for data-intensive models to

401     improve the estimates of population exposures to air pollution. Another important finding was

402     that the highest concentrations (10 to 20 $\mu g/m^3$) were estimated in the Brisbane, Gold Coast and

403     Logan City which have the highest population density in SEQ (Fig. 5); hence, the increased risk

404     of population expose to higher concentration of $PM_{2.5}$. Although the average $PM_{2.5}$ concentration

405     was below the WHO guideline (25 $\mu g/m^3$), epidemiological studies have demonstrated that $PM_{2.5}$

406     exposure to even lower concentrations is associated with high health risks (Burnett et al., 2014).

407     Based on the shape of the $PM_{2.5}$ exposure-response curves derived by Burnett et al. and used in

408     the Global Burden of Disease studies (Burnett et al., 2014), characterizing with a steep slope for

409     concentrations ranging from 10 to 20 $\mu g/m^3$, exposure to $PM_{2.5}$ concentration between 10 to 20

410     $\mu g/m^3$ highly increases the relative risk of stroke and chronic obstructive pulmonary disease;

411     therefore, even $PM_{2.5}$ concentration below the WHO guideline could not be considered safe.

412     Different methodologies make it difficult to compare our results to other studies, however we

413     have attempted to compare our results with two studies which have demonstrated the ability of

22

414  remotely sensed AOD and meteorological data to predict $PM_{2.5}$ concentration (Gupta and

415  Christopher, 2009a; Wu et al., 2012). Both studies used BPANN method to estimate the

416  spatiotemporal variation of $PM_{2.5}$, and reported $R^2$ lower than 0.61. Our model exhibited better

417  correlations than these models, which could be due to either: (1) the comprehensive input

418  variables used or (2) the more robust modeling algorithms used. We also compared our

419  methodology with a study of the global burden of disease 2013 conducted by Brauer et. al.

420  (2016) which combined ground measurements, chemical transport model outputs, and satellite-

421  based data to provide global estimates of $PM_{2.5}$ concentration. Although, chemical transport

422  model simulations were unavailable for our study area, our model was still able to capture 81%

423  of $PM_{2.5}$ variations.

424      Previous research demonstrated that PBLH and WS significantly affect the $PM_{2.5}$-AOD

425  relationship. Our results support these findings, but also demonstrate that incorporating other

426  spatial and spatiotemporal data as well as road density, land use types, NDVI, and industrial

427  point sources improves the model's performance.

428  **5. Conclusions**

429      Three different soft computing methods were applied to develop a satellite-based model for

430  estimating the spatiotemporal variation of $PM_{2.5}$. ANFIS performed very well compared to SVM

431  and BPANN. It exhibited satisfactory performance with CV-$R^2$, and CV-RMSE equals to 0.81,

432  and 1.79 $\mu g/m^3$, respectively. It provides estimates of monthly $PM_{2.5}$ concentrations during 2006-

433  2011. The modelling approach used in this study is highly applicable to similar settings

434  anywhere in the world assuming that researchers have access to data sets equivalent to these used

435  in our study. WRF, and its underlying boundary condition data, is a community model available

436 to the world research fraternity. The NASA Giovanni data is available to anyone in the research

437 community who can be registered with NASA as a data user. It is expected that researchers will

438 have access to all other similar data sets or proxies in their own national research and

439 information collection institutions, hence this method could be applied in other regions that

440 experience $PM_{2.5}$ exposure. We hope that our approach will be beneficial for epidemiological

441 studies and other researches seeking spatially accurate estimates of $PM_{2.5}$ with few monitoring

442 stations. It is certainly feasible to develop a model with higher spatial resolution which is a

443 direction of our future research. Further analysis such as global sensitivity and uncertainty

444 analyses (GSUA) can also be done to assess input factor importance and interaction (Lüdtke et

445 al., 2008; Saltelli et al., 2008). Recently developed Unified-Weather Research and Forecasting

446 model (NU-WRF) can be employed to obtain more accurate estimates of meteorological

447 parameters compared to WRF model (Peters-Lidard et al., 2015). Data management remains a

448 major challenge for empirical modelling as it requires to store, process and analyse large data

449 sets containing different types of data from multiple sources. Recently, new information models

450 have been developed to facilitate the data management and validation in observation-based

451 studies (Horsburgh et al., 2016; Shu et al., 2016).

**References**

463    Al-Alawi, S.M., Abdul-Wahab, S.A., Bakheit, C.S., 2008. Combining principal component
464    regression and artificial neural networks for more accurate predictions of ground-level ozone.
465    Environmental Modelling & Software 23(4) 396-403.
466    Australian Bureau of Statistics, 2011. Australian Statistical Geography Standard, Volume 1. ,
467    Main Structure and Greater Capital City Statistical Areas: Canbrra, Australia, p. 15.
468    Australian Bureau of Statistics, 2012. Year Book Australia: Canberra , Australia, pp. 237-251.
469    Australian Bureau of Statistics., 2012. Population Change in South-East Queensland: Australia.
470    Australian Government, 2015. Our Natural Environment: Australia.
471    Barnett, A.G., Williams, G.M., Schwartz, J., Neller, A.H., Best, T.L., Petroeschevsky, A.L.,
472    Simpson, R.W., 2005. Air pollution and child respiratory health: a case-crossover study in
473    Australia and New Zealand. American Journal of Respiratory and Critical Care Medicine
474    171(11) 1272-1278.
475    Beckerman, B.S., Jerrett, M., Martin, R.V., van Donkelaar, A., Ross, Z., Burnett, R.T., 2013a.
476    Application of the deletion/substitution/addition algorithm to selecting land use regression
477    models for interpolating air pollution measurements in California. Atmospheric Environment 77
478    172-177.
479    Beckerman, B.S., Jerrett, M., Serre, M., Martin, R.V., Lee, S.-J., van Donkelaar, A., Ross, Z., Su,
480    J., Burnett, R.T., 2013b. A Hybrid Approach to Estimating National Scale Spatiotemporal
481    Variability of PM2.5 in the Contiguous United States. Environmental Science & Technology
482    47(13) 7233-7241.
483    Beelen, R., Hoek, G., Vienneau, D., Eeftens, M., Dimakopoulou, K., Pedeli, X., Tsai, M.-Y.,
484    Künzli, N., Schikowski, T., Marcon, A., 2013. Development of $NO_2$ and NOx land use
485    regression models for estimating air pollution exposure in 36 study areas in Europe–the ESCAPE
486    project. Atmospheric Environment 72 10-23.
487    Brauer, M., Freedman, G., Frostad, J., van Donkelaar, A., Martin, R.V., Dentener, F., Dingenen,
488    R.v., Estep, K., Amini, H., Apte, J.S., Balakrishnan, K., Barregard, L., Broday, D., Feigin, V.,
489    Ghosh, S., Hopke, P.K., Knibbs, L.D., Kokubo, Y., Liu, Y., Ma, S., Morawska, L., Sangrador,
490    J.L.T., Shaddick, G., Anderson, H.R., Vos, T., Forouzanfar, M.H., Burnett, R.T., Cohen, A.,
491    2016. Ambient Air Pollution Exposure Estimation for the Global Burden of Disease 2013.
492    Environmental Science & Technology 50(1) 79-88.
493    Brook, R.D., Rajagopalan, S., Pope, C.A., Brook, J.R., Bhatnagar, A., Diez-Roux, A.V.,
494    Holguin, F., Hong, Y., Luepker, R.V., Mittleman, M.A., 2010. Particulate matter air pollution
495    and cardiovascular disease an update to the scientific statement from the American Heart
496    Association. Circulation 121(21) 2331-2378.
497    Bureau of Meteorology, 2015. Normalised Difference Vegetation Index: Australia.

498  Burnett, R.T., Pope, C.A., Ezzati, M., Olives, C., Lim, S.S., Mehta, S., Shin, H.H., Singh, G.,
499  Hubbell, B., Brauer, M., 2014. An integrated risk function for estimating the global burden of
500  disease attributable to ambient fine particulate matter exposure. Environ Health Perspect. 122(4)
501  397-403.
502  Carnevale, C., Finzi, G., Pederzoli, A., Turrini, E., Volta, M., 2016. Lazy Learning based
503  surrogate models for air quality planning. Environmental Modelling & Software 83 47-57.
504  Chan, C.-C., Chuang, K.-J., Chien, L.-C., Chen, W.-J., Chang, W.-T., 2006. Urban air pollution
505  and emergency admissions for cerebrovascular diseases in Taipei, Taiwan. European heart
506  journal 27(10) 1238-1244.
507  Chen, L., Verrall, K., Tong, S., 2006. Air particulate pollution due to bushfires and respiratory
508  hospital admissions in Brisbane, Australia. International journal of environmental health research
509  16(03) 181-191.
510  Chu, D.A., Kaufman, Y., Zibordi, G., Chern, J., Mao, J., Li, C., Holben, B., 2003. Global
511  monitoring of air pollution over land from the Earth Observing System Terra Moderate
512  Resolution Imaging Spectroradiometer (MODIS). Journal of Geophysical Research:
513  Atmospheres (1984–2012) 108(D21).
514  Chudnovsky, A.A., Koutrakis, P., Kloog, I., Melly, S., Nordio, F., Lyapustin, A., Wang, Y.,
515  Schwartz, J., 2014. Fine particulate matter predictions using high resolution Aerosol Optical
516  Depth (AOD) retrievals. Atmospheric Environment 89 189-198.
517  Chudnovsky, A.A., Lee, H.J., Kostinski, A., Kotlov, T., Koutrakis, P., 2012. Prediction of daily
518  fine particulate matter concentrations using aerosol optical depth retrievals from the
519  Geostationary Operational Environmental Satellite (GOES). Journal of the Air & Waste
520  Management Association 62(9) 1022-1031.
521  Corani, G., Scanagatta, M., 2016. Air pollution prediction via multi-label classification.
522  Environmental Modelling & Software 80 259-264.
523  Dirgawati, M., Barnes, R., Wheeler, A.J., Arnold, A.-L., McCaul, K.A., Stuart, A.L., Blake, D.,
524  Hinwood, A., Yeap, B.B., Heyworth, J.S., 2015. Development of Land Use Regression models
525  for predicting exposure to NO2 and NOx in Metropolitan Perth, Western Australia.
526  Environmental Modelling & Software 74 258-267.
527  Forouzanfar, M.H., Alexander, L., Anderson, H.R., Bachman, V.F., Biryukov, S., Brauer, M.,
528  Burnett, R., Casey, D., Coates, M.M., Cohen, A., 2015. Global, regional, and national
529  comparative risk assessment of 79 behavioural, environmental and occupational, and metabolic
530  risks or clusters of risks in 188 countries, 1990–2013: a systematic analysis for the Global
531  Burden of Disease Study 2013. The Lancet 386(10010) 2287-2323.
532  Franklin, M., Zeka, A., Schwartz, J., 2007. Association between $PM_{2.5}$ and all-cause and
533  specific-cause mortality in 27 US communities. Journal of Exposure Science and Environmental
534  Epidemiology 17(3) 279-287.
535  Giavarina, D., 2015. Understanding Bland Altman analysis. Biochemia medica 25(2) 141-151.
536  Gupta, P., Christopher, S.A., 2009a. Particulate matter air quality assessment using integrated
537  surface, satellite, and meteorological products: 2. A neural network approach. Journal of
538  Geophysical Research: Atmospheres (1984–2012) 114(D20).
539  Gupta, P., Christopher, S.A., 2009b. Particulate matter air quality assessment using integrated
540  surface, satellite, and meteorological products: Multiple regression approach. Journal of
541  Geophysical Research: Atmospheres (1984–2012) 114(D14).

542     Gupta, P., Christopher, S.A., Box, M.A., Box, G.P., 2007. Multi year satellite remote sensing of
543     particulate matter air quality over Sydney, Australia. International Journal of Remote Sensing
544     28(20) 4483-4498.
545     Gupta, P., Christopher, S.A., Wang, J., Gehrig, R., Lee, Y., Kumar, N., 2006. Satellite remote
546     sensing of particulate matter and air quality assessment over global cities. Atmospheric
547     Environment 40(30) 5880-5892.
548     Han, L., Zhou, W., Li, W., 2015. Increasing impact of urban fine particles (PM2.5) on areas
549     surrounding Chinese cities. Scientific reports 5 12467.
550     Harrison, R.M., Deacon, A.R., Jones, M.R., Appleby, R.S., 1997. Sources and processes
551     affecting concentrations of PM10 and PM2.5 particulate matter in Birmingham (UK).
552     Atmospheric Environment 31(24) 4103-4117.
553     Henderson, S.B., Beckerman, B., Jerrett, M., Brauer, M., 2007. Application of land use
554     regression to estimate long-term concentrations of tra☐c-related nitrogen oxides and fine
555     particulate matter. Environmental Science & Technology 41(7) 2422-2428.
556     Hinwood, A., De Klerk, N., Rodriguez, C., Jacoby, P., Runnion, T., Rye, P., Landau, L., Murray,
557     F., Feldwick, M., Spickett, J., 2006. The relationship between changes in daily air pollution and
558     hospitalizations in Perth, Australia 1992–1998: a case-crossover study. International journal of
559     environmental health research 16(1) 27-46.
560     Hoff, R.M., Christopher, S.A., 2009. Remote sensing of particulate pollution from space: have
561     we reached the promised land? Journal of the Air & Waste Management Association 59(6) 645-
562     675.
563     Hoffmann, B., Moebus, S., Möhlenkamp, S., Stang, A., Lehmann, N., Dragano, N.,
564     Schmermund, A., Memmesheimer, M., Mann, K., Erbel, R., 2007. Residential exposure to traffic
565     is associated with coronary atherosclerosis. Circulation 116(5) 489-496.
566     Horsburgh, J.S., Aufdenkampe, A.K., Mayorga, E., Lehnert, K.A., Hsu, L., Song, L., Jones, A.S.,
567     Damiano, S.G., Tarboton, D.G., Valentine, D., 2016. Observations Data Model 2: A community
568     information model for spatially discrete Earth observations. Environmental Modelling &
569     Software 79 55-74.
570     Hu, Z., Follmann, D.A., Miura, K., 2015. Vaccine design via nonnegative lasso☐based variable
571     selection. Statistics in medicine 34(10) 1791-1798.
572     Jang, J.-S.R., 1993. ANFIS: adaptive-network-based fuzzy inference system. Systems, Man and
573     Cybernetics, IEEE Transactions on 23(3) 665-685.
574     Jerrett, M., Arain, A., Kanaroglou, P., Beckerman, B., Potoglou, D., Sahsuvaroglu, T., Morrison,
575     J., Giovis, C., 2005a. A review and evaluation of intraurban air pollution exposure models.
576     Journal of Exposure Science and Environmental Epidemiology 15(2) 185-204.
577     Jerrett, M., Burnett, R.T., Ma, R., Pope, C.A., Krewski, D., Newbold, K.B., Thurston, G., Shi,
578     Y., Finkelstein, N., Calle, E.E., Thun, M.J., 2005b. Spatial analysis of air pollution and mortality
579     in Los Angeles. Epidemiology 16(6) 727-736.
580     Jones, D.A., Wang, W., Fawcett, R., 2009. High-quality spatial climate data-sets for Australia.
581     Australian Meteorological and Oceanographic Journal 58(4) 233.
582     Kaufman, Y.J., Tanré, D., Boucher, O., 2002. A satellite view of aerosols in the climate system.
583     Nature 419(6903) 215-223.
584     Kim, J.-H., 2009. Estimating classification error rate: Repeated cross-validation, repeated hold-
585     out and bootstrap. Computational Statistics & Data Analysis 53(11) 3735-3745.

586 Kloog, I., Koutrakis, P., Coull, B.A., Lee, H.J., Schwartz, J., 2011. Assessing temporally and
587 spatially resolved PM2.5 exposures for epidemiological studies using satellite aerosol optical
588 depth measurements. Atmospheric Environment 45(35) 6267-6275.
589 Knibbs, L.D., Hewson, M.G., Bechle, M.J., Marshall, J.D., Barnett, A.G., 2014. A national
590 satellite-based land-use regression model for air pollution exposure assessment in Australia.
591 Environmental research 135 204-211.
592 Kondragunta, S., Lee, P., McQueen, J., Kittaka, C., Prados, A.I., Ciren, P., Laszlo, I., Pierce,
593 R.B., Hoff, R., Szykman, J.J., 2008. Air quality forecast verification using satellite data. Journal
594 of Applied Meteorology and Climatology 47(2) 425-442.
595 Kruse, R., Pasi, G., Alonso, J.M., 2013. Introduction to the Soft Computing and Intelligent Data
596 Analysis Minitrack, 46th Hawaii International Conference on System Sciences, p. 1384.
597 Kunzli, N., Jerrett, M., Mack, W.J., Beckerman, B., LaBree, L., Gilliland, F., Thomas, D., Peters,
598 J., Hodis, H.N., 2005. Ambient air pollution and atherosclerosis in Los Angeles. Environ Health
599 Perspect. 113(2) 201-206.
600 Lee, J.-H., Wu, C.-F., Hoek, G., de Hoogh, K., Beelen, R., Brunekreef, B., Chan, C.-C., 2014.
601 Land use regression models for estimating individual NOx and NO2 exposures in a metropolis
602 with a high density of traffic roads and population. Science of The Total Environment 472 1163-
603 1171.
604 Levelt, P.F., Van den Oord, G.H., Dobber, M.R., Mälkki, A., Visser, H., De Vries, J., Stammes,
605 P., Lundell, J.O., Saari, H., 2006. The ozone monitoring instrument. Geoscience and Remote
606 Sensing, IEEE Transactions on 44(5) 1093-1101.
607 Li, C., Hsu, N.C., Tsay, S.-C., 2011. A study on the potential applications of satellite data in air
608 quality monitoring and forecasting. Atmospheric Environment 45(22) 3663-3675.
609 Li, Q., Shao, J., 2015. Regularizing lasso: a consistent variable selection method. Statistica
610 Sinica 25 975-992.
611 Lin, C.-T., Lee, C.S.G., 1991. Neural-network-based fuzzy logic control and decision system.
612 IEEE Transactions on computers 40(12) 1320-1336.
613 Lin, C., Li, Y., Yuan, Z., Lau, A.K.H., Li, C., Fung, J.C.H., 2014. Using satellite remote sensing
614 data to estimate the high-resolution distribution of ground-level PM2.5. Remote Sensing of
615 Environment 156 117-128.
616 Lindstrom, J., Szpiro, A.A., Sampson, P.D., Sheppard, L., Oron, A.P., Richards, M., Larson, T.,
617 2011. A Flexible Spatio-Temporal Model for Air Pollution: Allowing for Spatio-Temporal
618 Covariates. Berkeley Electronics Press.
619 Liu, Y., Paciorek, C.J., Koutrakis, P., 2009. Estimating regional spatial and temporal variability
620 of PM2.5 concentrations using satellite data, meteorology, and land use information.
621 Lüdtke, N., Panzeri, S., Brown, M., Broomhead, D.S., Knowles, J., Montemurro, M.A., Kell,
622 D.B., 2008. Information-theoretic sensitivity analysis: a general method for credit assignment in
623 complex networks. Journal of The Royal Society Interface 5(19) 223-235.
624 Lunn, H., Johnston, C., Flavel, R., 2002. The Vision and Living Skills Research Project: Meeting
625 the challenge of intervention in urban and rural communities, 11th ICEVI World Conference:
626 Noordwijkerhout, the Netherlands.
627 Maciejewska, K., Juda-Rezler, K., Reizer, M., Klejnowski, K., 2015. Modelling of black carbon
628 statistical distribution and return periods of extreme concentrations. Environmental Modelling &
629 Software 74 212-226.
630 McBride, S.J., Williams, R.W., Creason, J., 2007. Bayesian hierarchical modeling of personal
631 exposure to particulate matter. Atmospheric Environment 41(29) 6143-6155.

Meyer, C.M., Luhar, A.K., Mitchell, R.M., 2008. Biomass burning emissions over northern Australia constrained by aerosol measurements: I—Modelling the distribution of hourly emissions. Atmospheric Environment 42(7) 1629-1646.

Moazami, S., Noori, R., Amiri, B.J., Yeganeh, B., Partani, S., Safavi, S., 2016. Reliable prediction of carbon monoxide using developed support vector machine. Atmospheric Pollution Research 7(3) 412-418.

NASA, 2013. Global Change Master Directory, Earth Science Keywords, Aerosol and Warming, 8th ed: USA.

Novotny, E.V., Bechle, M.J., Millet, D.B., Marshall, J.D., 2011. National satellite-based land-use regression: NO2 in the United States. Environmental Science & Technology 45(10) 4407-4414.

Ordieres, J., Vergara, E., Capuz, R., Salazar, R., 2005. Neural network prediction model for fine particulate matter (PM2.5) on the US–Mexico border in El Paso (Texas) and Ciudad Juárez (Chihuahua). Environmental Modelling & Software 20(5) 547-559.

Ovaska, S.J., 2004. Computationally Intelligent Hybrid Systems: The Fusion of Soft Computing and Hard Computing (IEEE Press Series on Computational Intelligence). Wiley-IEEE Press.

Pandey, G., Zhang, B., Jian, L., 2013. Predicting submicron air pollution indicators: a machine learning approach. Environmental Science: Processes & Impacts 15(5) 996-1005.

Peters-Lidard, C.D., Kemp, E.M., Matsui, T., Santanello, J.A., Kumar, S.V., Jacob, J.P., Clune, T., Tao, W.-K., Chin, M., Hou, A., 2015. Integrated modeling of aerosol, cloud, precipitation and land processes at satellite-resolved scales. Environmental Modelling & Software 67 149-159.

Queensland Government., 2014. Air Quality Monitoring. Queensland Government: Australia.

Refaeilzadeh, P., Tang, L., Liu, H., 2009. Cross-validation, Encyclopedia of database systems. Springer, pp. 532-538.

Reid, C.E., Jerrett, M., Petersen, M.L., Pfister, G.G., Morefield, P.E., Tager, I.B., Raffuse, S.M., Balmes, J.R., 2015. Spatiotemporal Prediction of Fine Particulate Matter During the 2008 Northern California Wildfires Using Machine Learning. Environmental Science & Technology 49(6) 3887-3896.

Reis, S., Seto, E., Northcross, A., Quinn, N.W., Convertino, M., Jones, R.L., Maier, H.R., Schlink, U., Steinle, S., Vieno, M., 2015. Integrating modelling and smart sensors for environmental and human health. Environmental Modelling & Software 74 238-246.

Rose, N., Cowie, C., Gillett, R., Marks, G.B., 2010. Validation of a spatiotemporal land use regression model incorporating fixed site monitors. Environmental Science & Technology 45(1) 294-299.

Salam, M.T., Islam, T., Gilliland, F.D., 2008. Recent evidence for adverse effects of residential proximity to traffic sources on asthma. Current opinion in pulmonary medicine 14(1) 3-8.

Saltelli, A., Ratto, M., Andres, T., Campolongo, F., Cariboni, J., Gatelli, D., Saisana, M., Tarantola, S., 2008. Global sensitivity analysis: the primer. John Wiley & Sons.

Samet, J.M., 2007. Traffic, air pollution, and health. Inhalation toxicology 19(12) 1021-1027.

Shu, Y., Liu, Q., Taylor, K., 2016. Semantic validation of environmental observations data. Environmental Modelling & Software 79 10-21.

Simpson, R., Williams, G., Petroeschevsky, A., Best, T., Morgan, G., Denison, L., Hinwood, A., Neville, G., Neller, A., 2005. The short□term effects of air pollution on daily mortality in four Australian cities. Australian and New Zealand Journal of Public Health 29(3) 205-212.

Strawa, A., Chatfield, R., Legg, M., Scarnato, B., Esswein, R., 2013. Improving Retrievals of Regional PM2.5 Concentrations From MODIS and OMI Multi-Satellite Observations, AGU Fall Meeting Abstracts, p. 0300.

678    Su, J., Jerrett, M., Beckerman, B., 2009. A distance-decay variable selection strategy for land use
679    regression modeling of ambient air pollution exposures. Science of The Total Environment
680    407(12) 3890-3898.
681    Syed, A.R., 2011. A review of cross validation and adaptive model selection, Department of
682    Mathematics and Statistics. Georgia State University.
683    Tibshirani, R., 2011. Regression shrinkage and selection via the lasso: a retrospective. Journal of
684    the Royal Statistical Society: Series B (Statistical Methodology) 73(3) 273-282.
685    Torres, O., Tanskanen, A., Veihelmann, B., Ahn, C., Braak, R., Bhartia, P.K., Veefkind, P.,
686    Levelt, P., 2007. Aerosols and surface UV products from Ozone Monitoring Instrument
687    observations: An overview. Journal of Geophysical Research: Atmospheres (1984–2012)
688    112(D24).
689    United Nations, 2015. World Population Prospects: The 2015 Revision. Department of
690    Economic and Social Affairs, Population Division.
691    Vienneau, D., de Hoogh, K., Bechle, M.J., Beelen, R., van Donkelaar, A., Martin, R.V., Millet,
692    D.B., Hoek, G., Marshall, J.D., 2013. Western European land use regression incorporating
693    satellite-and ground-based measurements of NO2 and PM10. Environmental Science &
694    Technology 47(23) 13555-13564.
695    Wu, Y., Guo, J., Zhang, X., Tian, X., Zhang, J., Wang, Y., Duan, J., Li, X., 2012. Synergy of
696    satellite and ground based observations in estimation of particulate matter in eastern China.
697    Science of The Total Environment 433 20-30.
698    Yeganeh, B., Motlagh, M.S.P., Rashidi, Y., Kamalan, H., 2012. Prediction of CO concentrations
699    based on a hybrid partial least square and support vector machine model. Atmospheric
700    Environment 55 357-365.

701