

۱ **Satellite-based prediction of surface dust mass concentration in southeastern**
۲ **Iran using gradient boosting regression**

۳
۴ Seyed Babak Haji Seyed Asadollah¹, Ahmad Sharafati^{1*}, Davide Motta², Antonio Jodar-Abellan^{3,4},
۵ Miguel Ángel Pardo^{4,5}

۶ ¹ Department of Civil Engineering, Science and Research Branch, Islamic Azad University, Tehran, Iran.

۷ ² Department of Mechanical and Construction Engineering, Northumbria University, Wynne Jones
۸ Building, Newcastle upon Tyne NE1 8ST, United Kingdom

۹ ³ Departamento de Análisis Geográfico Regional y Geografía Física, Facultad de Filosofía y Letras,
۱۰ Campus Universitario de Cartuja, University of Granada, 18071 Granada, Spain

۱۱ ⁴ University institute of Water and Environmental Sciences. University of Alicante, Spain

۱۲ ⁵ Department of Civil Engineering. University of Alicante, Spain

۱۳ **Corresponding authors:** Ahmad Sharafati

۱۴ **Corresponding email:** asharafati@srbiau.ac.ir, asharafati@gmail.com

۲۴ **Abstract**

۲۵ The southeastern section of Iran, especially the province of Khuzestan, experience
۲۶ severe air pollution levels, such as high values of Surface Dust Mass Concentration
۲۷ (SDMC). The province lacks accurate and well-placed ground observational
۲۸ stations, therefore the only viable approach for evaluating SDMC is via remote
۲۹ sensing. In this study, meteorological, hydrological and geological data on 11 input
۳۰ variables from Modern-Era Retrospective analysis for Research and Applications
۳۱ Version (MERRA-2), Global Precipitation Measurement (GPM) and Global Land
۳۲ Data Assimilation System (GLDAS) for the year 2018 are used for prediction of the
۳۳ SDMC values, also obtained from MERRA-2. For real-time prediction, Pearson's
۳۴ Correlation Coefficient (PCC) analysis shows that wind-related variables – surface
۳۵ wind speed, surface aerodynamic conductivity and surface pressure – are those with
۳۶ the highest correlation with SDMC. Using the Gradient Boosting Regression (GBR)
۳۷ algorithm, these three variables can simulate SDMC with good accuracy ($CC =$
۳۸ 0.815 , $N - RMSE = 0.605$). Future forecasting of SDMC requires knowledge of
۳۹ both wind-related and heat-related variables. However, SDMC predictions can be
۴۰ obtained with the GBR algorithm with adequate accuracy ($CC = 0.640$, $NRMSE =$
۴۱ 0.781) by just considering the surface pressure value observed four days before the
۴۲ forecasted day. This study shows that robust predictions of SDMC can be obtained
۴۳ using exclusively remote sensing data, without ground-based observations.

۴۴

۴۵

۴۶ **Keywords:** Surface dust mass concentration, Gradient Boosting Regression,
۴۷ Prediction, Khuzestan

۴۸

۴۹

۵۰

1. Introduction

Recent reports by the World Health Organization (WHO) revealed that nearly seven million people die in the world due to high levels of air pollution every year (WHO, 2014). Major air pollutants are carbon monoxide, sulfur, nitrogen dioxide, and surface-level ozone (Chen *et al.*, 2007, Council, 1992, Duan *et al.*, 2019, Sunyer *et al.*, 2003). Particulate Matter (PM) is also highlighted by WHO as the thirteenth mortality cause around the world, which makes it further hazardous (Anderson *et al.*, 2012). PM is a mixture of microscopic particles and liquid substances such as metals, organic materials, acids and dust (Planning, 1996, USEPA, 2019). PM can cause significant health problems, such as blood pressure, lung cancer, and cardiovascular diseases (Brook & Rajagopalan, 2009, Hamanaka & Mutlu, 2018, Raaschou-Nielsen *et al.*, 2016), therefore monitoring its concentration is critical for health and environmental purposes.

Numerous studies have been conducted on PM concentration measurement, using various sampling approaches and instruments (Amaral *et al.*, 2015, Kwasny *et al.*, 2010, Nakata *et al.*, 2013). Ground-based measurement, mostly done with monitoring stations, is the most commonly employed method. However, while the direct sampling approach is considered as the most accurate, it is neither cost- nor time-efficient, and station numbers are typically limited and station spatial distribution is generally irregular (M. Lee *et al.*, 2016, D. Liu & Li, 2015).

Recently, to overcome the disadvantages of direct sampling, Remote Sensing (RS) methods have become increasingly utilized for forecasting of hydrological and meteorological phenomena (Asadollah *et al.*, 2021, Ghizat *et al.*, 2022, Shiru *et al.*, 2022). RS covers a wide range of spatial and temporal data observations that can particularly benefit regions that lack direct observation stations (Campbell & Wynne, 2011, Davis & Swain, 1978), as in the case of Iran. To the authors'

۷۷ knowledge, there are no useful historical ground-based datasets regarding the
۷۸ concentration of air pollutants in any province of Iran.

۷۹ Early satellite versions were unable to record near-surface PM concentration and
۸۰ instead provided a substitute parameter called Aerosol Optical Depth or AOD (Diao
۸۱ *et al.*, 2019). Considering a ray of light being radiated from a satellite source, the
۸۲ AOD is defined as the decay level of that light reflection from the surface, which is
۸۳ mainly caused by the presence of particles in the air column (Van Donkelaar *et al.*,
۸۴ 2010). The Multi-Angle Implementation of Atmospheric Correction (MAIAC)
۸۵ method, used in conjunction with the satellite-based sensor Moderate Resolution
۸۶ Imaging Spectroradiometer (MODIS), is widely used for predicting the AOD over
۸۷ different regions around the globe. (A. Chudnovsky *et al.*, 2013) evaluated the
۸۸ applicability of the MAIAC algorithm by comparing its predictions with ground-
۸۹ based PM_{2.5} observations from 84 monitoring stations across New England in the
۹۰ United States over the period 2002-2008. The results indicated that the AOD
۹۱ obtained from the MAIAC is correlated with the observed PM_{2.5} surface
۹۲ concentration. This research also showed that the MAIAC-AOD shows a better
۹۳ correlation with the in-situ PM_{2.5} compared to conventional MODIS-AQUA
۹۴ products.

۹۵ Several studies focused on finding the relation between the AOD and ground-level
۹۶ PM concentration. (H. J. Lee *et al.*, 2012) developed a statistical method to predict
۹۷ the concentration of daily PM_{2.5} by combining the satellite AOD data with the
۹۸ ground-based ones. Specifically, the MODIS satellite outputs were used together
۹۹ with the observed data from the U.S. Environmental Protection Agency (EPA). The
۱۰۰ authors evaluated two groups of days consisting of days with or without satellite data
۱۰۱ availability during the period 2000-2008 for the New England. With a Pearson's
۱۰۲ Correlation Coefficient (PCC) of 0.91, their predictions were showed this to be a

1.03 suitable approach, especially in urban areas. (A. A. Chudnovsky *et al.*, 2014)
1.04 predicted the Fine Particulate Matter (FPM) in the air using high-resolution aerosol
1.05 data obtained with the MAIAC algorithm using MODIS satellite observations.
1.06 Several meteorological parameters (e.g., speed of wind and relative humidity), as
1.07 well as the land use, were utilized to predict the daily-based FPM over New England.
1.08 After calibrating the FPM satellite data with ground-based observation originated
1.09 from the EPA, they used a novel interpolation approach, the Inverse Probability
1.10 Weighting (IPW), to complete the prediction task. Their proposed model predicted
1.11 real-time FPM with high accuracy. Analogously, (Just *et al.*, 2015) used the MODIS
1.12 daily AOD values to predict PM_{2.5} over Mexico City from 2004 to 2014 using
1.13 statistical modeling. With a correlation coefficient R value of 0.85, their model
1.14 proved to be an accurate tool for predicting PM_{2.5} concentration an subcategory of
1.15 AOD. (X. Zhang *et al.*, 2018) employed the MAIAC-MODIS satellite outputs to
1.16 extract the AOD records and developed a multi-input statistical model based on
1.17 geographical properties, climate variables (air temperature, wind speed, and
1.18 visibility), and land use data to predict the ground-measured PM_{2.5} concentrations
1.19 over Texas in the United States, between the years 2008 and 2013. Their proposed
1.20 model provided accurate predictions with a correlation coefficient of 0.79~0.83.

1.21 In the last decade, technological advancement has led Artificial Intelligence (AI) to
1.22 become the dominant regression and classification approach in many research fields
1.23 (Mehdizadeh *et al.*, 2017, Nourani *et al.*, 2014, W.-C. Wang *et al.*, 2009). Compared
1.24 to statistical and numerical methods, AI can achieve the desired target in a much
1.25 faster and easier manner (Al-Othman *et al.*, 2022, Karandish & Šimůnek, 2016).
1.26 Benchmark AI algorithms such as Artificial Neural Network (ANN) and Adaptive
1.27 Neuro-Fuzzy Inference System (ANFIS) have been widely used for prediction in
1.28 earth sciences. (Mirzaei *et al.*, 2019) investigated the relationship between the

129 satellite-originated AOD values and ground measured PM2.5 concentrations over
130 Tehran in Iran. A model known as Geographically and Temporally Weighted
131 Regression (GTWR) was used to assess this relationship between the years of 2011
132 and 2017 and convert MODIS-AOD values to PM2.5 surface concentrations.
133 Comparison of four different AI algorithms reveal that the Generalized Regression
134 Neural Network (GRNN) algorithm performed better than its alternatives, ANN and
135 ANFIS.

136 More advanced AI methods, based on Machine Learning (ML), generally provide
137 better prediction performance compared to “classic” AI algorithms. The MLs show
138 better task in reducing the prediction associated bias and variances, have better
139 overfitting-elusive procedures and can be manually tuned more easily (Khazode &
140 Sarode, 2020, Müller & Guido, 2016, Wuest *et al.*, 2016). While the early ML
141 models, such as Support Vector Machine (SVM) and Multivariate Adaptive
142 Regression Splines (MARS) demonstrated acceptable performance, newer models
143 called ensemble algorithms have shown superior applicability. Ensemble algorithms
144 such as Ada-boost, Random Forest (RF), and Extreme Tree Regression (ETR) were
145 successfully applied in various studies, outperforming the classic AI algorithms (F.
146 Wang *et al.*, 2021, J. Zhang *et al.*, 2019, Zhu *et al.*, 2021).

147 The literature review by (Chu *et al.*, 2016) shows that multiple studies have adopted
148 AI algorithms to predict aerosol levels over different regions of the world. For
149 example, (Di *et al.*, 2016) predicted the AOD over the Unites States using ANN
150 algorithms. (Nguyen *et al.*, 2015) employed Support Vector Regression (SVR) and
151 Multiple linear regression (MLR) to simulate the organic carbon concentrations in
152 Gosan, South Korea, between the years 2011 and 2012. Another related study (Lary
153 *et al.*, 2014) utilized a machine learning regression to predict the worldwide aerosol
154 concentration from 1997 to 2014. (Nabavi *et al.*, 2018) compared the performance

100 of several ML algorithms for the prediction of monthly AOD in the western region
106 of Asia using the MODIS outputs. They used wind characteristics, soil temperature,
107 rainfall, drought index, and several other parameters as ML initial predictors, while
108 the MODIS AOD value was used as the target variable. (Kianian *et al.*, 2021)
109 employed RF as an ensemble ML algorithm to predict the spatial distribution of
160 PM2.5, especially in regions that are prone to gaps in AOD coverage. They also used
161 a statistical approach known as Lattice Kriging. Like many previous studies, they
162 first calibrated the MODIS outputs with the EPA ground-based observations.
163 Surface pressure, wind components, temperature, rainfall, relative humidity,
164 radiation flux and many other variables were investigated as the meteorological
165 parameters for PM2.5 prediction.

166 While the majority of the previous studies used the AOD data obtained from the
167 MODIS satellite, few have investigated other satellite-based aerosol diagnosis
168 outputs such as those from the Modern-Era Retrospective analysis for Research and
169 Applications Version 2, known as MERRA-2 (Gelaro *et al.*, 2017). (Sun *et al.*, 2019)
170 compared the AOD outputs from MERRA-2 and MODIS and evaluated them
171 against ground-based observations at 12 stations in China. Their findings suggested
172 that the MERRA-2 outputs are in good agreement with both MODIS-based and
173 ground-based values. (Gueymard & Yang, 2020) focused on global AOD data for a
174 period of 15 years and showed that MERRA-2 performs better than the European
175 Centre for Medium-Range Weather Forecasts (ECMWF)'s Copernicus Atmosphere
176 Monitoring Service (CAMS). Besides AOD diagnosis outputs, MERRA-2 provides
177 data on ambient air pollutants such as sulfate and dust concentration at the surface
178 level.

179 This study focuses on surface dust, which is one of the main PM components, and
180 aims to forecast the Surface Dust Mass Concentration (SDMC). The dust-related

۱۸۱ output of MERRA-2 was selected as the target parameter in this study, as done in
۱۸۲ several other studies too (Ukhov *et al.*, 2020, Veselovskii *et al.*, 2018, Xu *et al.*,
۱۸۳ 2020, Yao *et al.*, 2020). The Gradient Boosting Regression (GBR) was used here,
۱۸۴ because it is a robust and effective ensemble-based prediction algorithm (Johnson *et*
۱۸۵ *al.*, 2018, Srivastava *et al.*, 2018, Y. Zhang & Haghani, 2015). Differently from
۱۸۶ previous investigations, this study uses exclusively remote sensing data for dust
۱۸۷ concentration prediction. This study also presents, for the first time, a comprehensive
۱۸۸ analysis of correlation between various meteorological, hydrological, and geological
۱۸۹ variables with SDMC, and successfully forecasts SDMC few days in advance.

۱۹۰

۱۹۱ **2. Materials and Methods**

۱۹۲ **2.1. Study Area**

۱۹۳ Dust bowl-like storms have a significant socio-environmental impact in Iran (Salami
۱۹۴ *et al.*, 2021). Several studies investigated short- and long-term AOD patterns over
۱۹۵ Iran (Arkian & Nicholson, 2018, Sabetghadam *et al.*, 2018, Salami *et al.*, 2021,
۱۹۶ Yousefi *et al.*, 2020). Nearly all these studies pinpoint the province of Khuzestan as
۱۹۷ the region with the highest dust concentration. For example, (Rezaei *et al.*, 2019)
۱۹۸ evaluated Iran based on its spatial and temporal dust aerosol patterns utilizing the
۱۹۹ MODIS outputs between 2006 to 2015. Their results show that the Khuzestan and
۲۰۰ Sistan provinces are the most affected provinces among others. Similar results were
۲۰۱ obtained by (Mirakbari & Ebrahimi Khusfi, 2020) and (Dadashi-Roudbari &
۲۰۲ Ahmadi, 2020), which makes the Khuzestan province a good case study for
۲۰۳ evaluating the SDMC.

۲۰۴ As shown in Figure 1, the Khuzestan province is located in the southwestern region
۲۰۵ of Iran. It has an approximate area of 63,000 km² and 4 million inhabitants. Based

۲۰۶ on the digital elevation map in Figure 1, Khuzestan’s elevation below and above the
۲۰۷ sea level ranges between -105 and 3741 meters, respectively. This is associated with
۲۰۸ great diversity in climate conditions, from the cold temperatures in the north to
۲۰۹ tropical conditions in the south. Khuzestan’s summer months are considered those
۲۱۰ from April to September, while October to March are the winter months. The annual
۲۱۱ average maximum and minimum temperature of this province are 50°C and 9°C in
۲۱۲ July and March, respectively. The annual precipitation rate varies from ~200 mm
۲۱۳ (sea coast in the south) to ~1050 mm (near the Zagros mountains in the north). The
۲۱۴ dominant wind direction in this province is from west to east and northwest to
۲۱۵ southeast (Zarasvandi *et al.*, 2011).

[Figure 1]

۲۱۷ Dust storms, affecting air quality and impacting social life and economy, have
۲۱۸ become increasingly frequent in Khuzestan. These storms mainly happen in the
۲۱۹ summer months and originate from neighboring countries such as Iraq, with west
۲۲۰ winds into Iran (Daniali & Karimi, 2019). Based on reports from the Khuzestan
۲۲۱ meteorological stations, just in the year 2008 total of 1035 dust storm events were
۲۲۲ reported, which is considered a significant number (Zarasvandi, 2009).

۲۲۴ 2.2. Satellite Data

۲۲۵ 2.2.1. Modern-Era Retrospective Analysis for Research and Applications ۲۲۶ Version 2 (MERRA-2)

۲۲۷ Due to the high improvement in assimilation structure, The MERRA-2 replaced the
۲۲۸ original MERRA. It has a more advanced system including hyperspectral radiance
۲۲۹ and microwave examination. It also includes the Goddard Earth Observing System
۲۳۰ (GEOS-5) upgrade and ozone samplings of NASA, which makes it an applicable

۲۳۱ climate evaluation tool (Gelaro *et al.*, 2017). MERRA-2 benefits from the
۲۳۲ employment of a Grid-point Statistical Interpolation (GSI) climate analysis program,
۲۳۳ which is structured based on an additive analysis procedure that evaluate the
۲۳۴ incremental the meteorological data every 6 hours (Gelaro *et al.*, 2017, MERRA,
۲۳۵ 2AD).

۲۳۶ In this study, several types of MERRA-2 outputs were used, from different
۲۳۷ databases. First, the aerosol diagnosis from M2T1NXAER was used to extract
۲۳۸ surface dust mass concentration ($\frac{mg}{m^3}$) values over Iran. The M2T1NXAER has a
۲۳۹ temporal resolution of 1 hour and a spatial resolution is $0.5^\circ \times 0.625^\circ$ longitude and
۲۴۰ latitude, respectively. M2T1NXAER records black carbon, dry dust, organic carbon,
۲۴۱ sea salt, and sulfate aerosols in the air (Randles *et al.*, 2017). M2I1NXLFO, also
۲۴۲ used in this paper, mainly includes land surface data. It has a Same spatial and
۲۴۳ temporal resolution to M2T1NXAER and includes parameters such as surface layer
۲۴۴ height, pressure, air temperature, wind speed, and specific humidity (Reichle *et al.*,
۲۴۵ 2017).

۲۴۶ **2.2.2. Global Land Data Assimilation System (GLDAS)**

۲۴۷ GLDAS Version 2, used in this study, is structured in three components, GLDAS-
۲۴۸ 2.0, -2.1, and -2.2. The former, GLDAS-2.0 is completely in congruity with the
۲۴۹ Princeton meteorological observations and covers the period 1948 to 2014. The year
۲۵۰ 2000 to present is covered by the 2.1 version. Unlike the two mentioned versions,
۲۵۱ the GLDAS-2.2 observations utilized data adjustment. GLDAS-2.1 has two major
۲۵۲ streams, one is associated with the Global Precipitation Climatology Project (GPCP)
۲۵۳ precipitation products, and one is operating without it. The reason behind this is the
۲۵۴ 3 to 4-months postponement of GPCP, which forced version 2.1 to represent a
۲۵۵ temporary data without it called early products. Once the GPCP product become

256 accessible the GLDAS become synchronize with it and the early products become
 257 as archive. The data used in this study are from GLDAS-2.1 and have a temporal
 258 resolution of 3 hours and a spatial resolution of 0.25° . This product is simulated
 259 with version 7 of Land Information System (LIS) from model 3.6 of NOAH. In late
 260 2020 the 3-hourly and monthly GLDAS-2.0 products were re processed with the
 261 land mask data from MODIS-MOD44W which corrected the previous version issues
 262 such as data missing (Rodell *et al.*, 2004).

263 **2.3. Gradient Boosting Regression (GBR) Algorithm**

264 The boosting technique is based on aggregating a set of simple predictors. This
 265 aggregation procedure is structured by focusing on errors originated in each step till
 266 a better predictor is constructed with the least outcome error (Nie *et al.*, 2021).
 267 Considering Y as the target variable and $X = \{X_1, X_2, \dots, X_n\}$ as the input variables,
 268 the aim of the algorithm is to approximate $G'(X)$ as a branch of the original function
 269 $G(X)$ to map X to Y , so that the loss function $\mathcal{L}(Y, G(X))$ becomes minimum.

$$G'(X) = \operatorname{argmin} \mathcal{L}_{Y,X}(Y, G(X)) \quad (1)$$

270 By fitting the simple predictors to the \mathcal{L} at each step of the regression procedure, the
 271 Gradient Boosting (GB) algorithm tries to reduce the errors characterizing the
 272 preceding steps. This error correctional strategy increases the prediction accuracy
 273 and simultaneously decreases the bias of the prediction model. Acknowledging the
 274 m ($m = 0, \dots, M$) as the number of Stages which GB takes to properly train a tree,
 275 algorithm first employs an initial simple predictor $G_{m=0}(X)$ and then enforces a
 276 gradient so that the \mathcal{L} is minimized. This gradient is computed as follows:

$$Y'_i = - \left[\frac{\partial \mathcal{L}(Y_i, G(X_i))}{\partial G(X_i)} \right]_{G(X)=G_{m-1}(X)} \quad (2)$$

277 Consider $T(X_i, \alpha)$ as a regression tree and α as the simple predictor, a new tree can
 278 be structured by solving Equation (3), where α_m is the simple predictor parameter
 279 at each stage and Ψ is their corresponding weight:

$$\alpha_m = \underset{\alpha}{\operatorname{argmin}} \sum_{i=1}^N [Y_i' - \Psi \times T(X_i, \alpha)]^2 \quad (3)$$

280
 281 Considering β_m as each stage's optimal length, the $G_m(X)$ is updated at each
 282 iteration m as follows:

$$\beta_m = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \mathcal{L}(Y_i - G_{m-1}(X_i) + \beta T(X_i, \alpha_m)) \quad (4)$$

$$G_m(X) = G_{m-1}(X) + \beta_m T(X_i, \alpha) \quad (5)$$

283
 284 The iteration continues until the $\beta_m T(X_i, \alpha)$ term in Equation (5) becomes its
 285 minimum possible value (Friedman, 2001). Figure shows a flowchart of the GBR
 286 algorithm.

287 **[Figure 2]**

288 In this study the GB algorithm was implemented using the ensemble sub-category
 289 of the Scikit-learn library of the Python Programming language. The provided
 290 Gradient Boosting Regression (GBR) algorithm has several parameters that need to
 291 be “tuned” so that the algorithm performs as its maximum accuracy. Table 1 shows
 292 the default and optimal values for each parameter, obtained with an iterative process
 293 developed by the authors.

294

Table 1: Gradient Boosting Regression (GBR) parameters.

Parameter	Description	Default Value	Optimal Value
loss	Enhancement approach for the loss function	Squared error	Absolute error
learning_rate	This parameter reduces the value of each tree's contribution	0.1	1.1
n_estimators	Number of boosting phases that must be accomplished	100	498
random_state	At each boosting step, this parameter sets the random seed delivered to each tree	None	34
max_depth	Individual regression estimators' maximum depth	3	1
min_samples_leaf	Number of minimum samples needed to consider a node a leaf node.	1	45

۲۹۶

۲۹۷ 3. Results

۲۹۸ 3.1. Selection of Year and Region of Focus

۲۹۹ MERRA-2-M2TMNXAER monthly SDMC data for Iran were downloaded from the
 ۳۰۰ NASA website for the period between 2009 and 2020
 ۳۰۱ (https://disc.gsfc.nasa.gov/datasets/M2TMNXAER_5.12.4/summary?keywords=M2TMNXAER). This
 ۳۰۲ satellite provides monthly averaged aerosol diagnosis with spatial resolution of
 ۳۰۳ $0.5^\circ \times 0.625^\circ$. Considering this resolution there were total of 1056 observational
 ۳۰۴ grid cross-section points over the Iran. By calculating the maximum value of SDMC
 ۳۰۵ among the 12 months in each specific year and at each of 1056 points, Figure 3
 ۳۰۶ shows the distribution of yearly maximum SDMC (mg/m^3) over Iran for the period

307 considered. For better understanding of each specific year ranking based on SDMC
308 level, thier Annual Maximum (A.M.) has been noted in the figure.

309 **[Figure 3]**

310 Figure 3 shows that the southern and southeastern sections of Iran, including the
311 provinces of Khuzestan and Hormozgan as well as Sistan and Baluchestan, are
312 highly prone to large dust concentrations. The Khuzestan province appears to be the
313 most impacted province and is the focus of this investigation. Over the 12-year
314 period considered, Iran was overall characterized by minimum, mean and maximum
315 SDMC values of 0.04, 0.295 and $1.36 \text{ mg}/\text{m}^3$ respectively. The year 2018, as the
316 year with highest SDMC observed value ($\sim 1.36 \text{ mg}/\text{m}^3$) was selected as the study
317 period in current study.

318 **3.2. Input Variable Correlation Analysis**

319 Having selected the year and region of focus, SDMC data for the year 2018 for the
320 Khuzestan province was extracted based on outputs extracted from MERRA-2-
321 M2T1NXADG

322 (https://disc.gsfc.nasa.gov/datasets/M2T1NXADG_5.12.4/summary?keywords=M2T1NXADG),

323 consisting of 1-hour time averaged aerosol diagnosis data with the same spatial
324 resolution as MERRA-2-M2TMNXAER. 1-hour data was converted to daily data
325 by taking the maximum value within each 24-hour period, therefore obtaining 365
326 SDMC values.

327 As mentioned, the Khuzestan province is characterized by different types of climate
328 conditions due to its significant variation in land elevation; therefore 18 Research
329 Points (RPs) were considered in this study, equally distributed across the province
330 (with latitudinal and longitudinal spacing of 0.5° and 0.625° , respectively, same as

the resolution of MERRA2-M2T1NXADG). These include the cold region in the north as well as the extremely hot regions in the south.

Then, several hydrological, meteorological and geological outputs from Global Precipitation Measurement (GPM), Global Land Data Assimilation System (GLDAS), MERRA-2-M2I1NXLFO and MERRA-2-M2T3NVASM with various spatial and temporal distribution were extracted. Table 2 lists these variables, which are the 11 input variables considered for SDMC prediction in this study, with their units and corresponding satellite.

Table 2: Input variables considered for SDMC prediction.

<i>Satellite</i>	<i>Variable</i>	<i>Description</i>	<i>Units</i>
<i>GLDAS</i>	<i>Acond</i>	Aerodynamic conductance	<i>m/s</i>
	<i>Esoil</i>	Evaporation flux from soil	<i>kg/m²s</i>
	<i>Qh</i>	Surface upward sensible heat flux	<i>W/m³</i>
	<i>Evap</i>	Evapotranspiration	<i>kg/m²s</i>
	<i>SoilMoist</i>	Surface soil moisture	<i>kg/m²</i>
<i>MERRA-2</i>	<i>PS</i>	Surface pressure	<i>Pa</i>
	<i>TLML</i>	Surface air temperature	<i>K</i>
	<i>SPEEDLML</i>	Surface wind speed	<i>m/s</i>
	<i>QLML</i>	Surface specific humidity	—
	<i>HLML</i>	Surface layer height	<i>m</i>
<i>GPM</i>	<i>Precip</i>	Precipitation rate	<i>mm</i>

The Pearson's Correlation Coefficient (PCC) was computed to quantify the correlation of each of the 11 variables in Table 2 with daily SDMC. While there are other alternatives, such as Spearman's or Kendall's rank correlation, various studies found the efficiency and robustness of PCC, especially in handling data with non-

345 apparent outliers and non-linearity (Chok, 2010, Hauke & Kossowski, 2011,
346 Rebekić *et al.*, 2015).

347 Figure 4 contains the relative PCC heat map for the 18 RPs, including the average
348 PCC values over all the RPs. The three parameters with the highest PCC value,
349 surface wind speed SPEEDLML ($PCC = 0.61$), aerodynamic conductance ACOND
350 ($PCC = 0.57$) and surface pressure PS ($PCC = 0.53$), were selected as input
351 variables for real-time daily SDMC prediction with the GBR algorithm. PS denotes
352 the atmospheric surface pressure that directly controls the movement of air masses
353 from regions with low pressure to regions with higher pressure (Gomis *et al.*, 2008,
354 Guo *et al.*, 2011). The surface aerodynamic conductance (ACOND) describes the
355 effect of surface roughness on the movement of air masses (S. Liu *et al.*, 2007,
356 Mallick *et al.*, 2018).

357 **[Figure 4]**

358 From Figure 4, the hydrological parameters rainfall, relative humidity and
359 evapotranspiration have the least correlation with SDMC in the Khuzestan province.

360 **3.3. Real-Time Prediction of SDMC**

361 The three input variables for prediction, known at the current time (day) “t”, were
362 used in the GBR algorithm to predict the current time SDMC. This was done for all
363 18 RPs and the prediction performance was evaluated using four indices, PCC,
364 Nash-Sutcliffe Efficiency (NSE), Normalized-Root Mean Squared Error (N-RMSE)
365 and Normalized Mean Absolute Error (N-MAE). Results of these metrics which
366 have been widely used in earth and water research fields (Jodar-Abellan *et al.*, 2019,
367 Moriasi *et al.*, 2007, Pardo *et al.*, 2020), are shown in Table 3.

368

۳۶۹ **Table 3:** SDMC real-time prediction performance for the 18 Research Points
 ۳۷۰ considered.

RP	PCC	NSE	N-RMSE	N-MAE
P1	0.716	0.493	0.710	0.494
P2	0.771	0.585	0.642	0.453
P3	0.815	0.632	0.605	0.422
P4	0.648	0.397	0.774	0.446
P5	0.707	0.473	0.723	0.516
P6	0.793	0.621	0.613	0.437
P7	0.783	0.598	0.632	0.455
P8	0.706	0.457	0.735	0.518
P9	0.720	0.490	0.712	0.507
P10	0.686	0.451	0.738	0.517
P11	0.505	0.170	0.908	0.613
P12	0.701	0.481	0.718	0.494
P13	0.684	0.425	0.756	0.563
P14	0.525	0.223	0.878	0.628
P15	0.405	0.079	0.957	0.594
P16	0.541	0.253	0.862	0.618
P17	0.435	0.139	0.925	0.702
P18	0.398	0.095	0.948	0.649

۳۷۱
 ۳۷۲ Values in Table 3 vary across the RPs but, overall, the average PCC, NSE, N-RMSE
 ۳۷۳ and N-MAE are within an acceptable range. The P3 Research Point ($CC = 0.815$,
 ۳۷۴ $NSE = 0.632$, $N - RMSE = 0.605$ and $N - MAE = 0.422$) is characterized by
 ۳۷۵ the best prediction performance. Figure 5 highlights a clear pattern of better
 ۳۷۶ prediction performance in the southern regions of the province. This figure has been
 ۳۷۷ obtained using the Inverse Distance Weighting (IDW) interpolation method over the
 ۳۷۸ Khuzestan province based on the calculated performance indices of 18 research
 ۳۷۹ points.

۳۸۰ **[Figure 5]**

۳۸۱

۳۸۲

۳۸۳ **3.4. Future Forecasting of SDMC**

۳۸۴ For the P3 Research Point (best real-time SDMC prediction, see previous section),
۳۸۵ future forecasting of SDMC was considered for lead times of ‘t-2’, ‘t-4’, ‘t-6’ and
۳۸۶ ‘t-8’ (input data from 2, 4, 6 and 8 days prior to the current time ‘t’ for which SDMC
۳۸۷ is forecasted). PCC values were calculated to quantify the correlation between the
۳۸۸ current time SDMC (at time “t”) and the 11 input variables from Table 2 for the four
۳۸۹ lead times. Figure 6 shows the PCC values in a circular bar chart.

۳۹۰ **[Figure 6]**

۳۹۱ From Figure 6 it can be seen that, analogously to the case of real-time prediction,
۳۹۲ precipitation and surface specific humidity have the lowest correlation with the
۳۹۳ SDMC. PS and surface upward sensible heat flux Qh and surface air temperature
۳۹۴ TLML correlations with SDMC show a significant increase when moving backward
۳۹۵ in time.

۳۹۶ Based on the PCC analysis and using a procedure based on progressive elimination
۳۹۷ of input variables with lower and lower PCC value, a procedure applied in other
۳۹۸ previous studies (Sharafati, Asadollah, & Hosseinzadeh, 2020, Sharafati, Asadollah,
۳۹۹ & Neshat, 2020), combinations of the 11 input variables from Table 2 were
۴۰۰ constructed for future forecasting of SDMC (lead times of ‘t-2’, ‘t-4’, ‘t-6’ and ‘t-
۴۰۱ 8’) using the GBR algorithm. The combinations considered are listed in Table 4.

۴۰۲

۴۰۳

۴۰۴

۴۰۵

۴۰۶

4.7

Table 4: Input variable combinations for future forecasting of SDMC.

Combination	Input Variables
C1	$PS(t-4), PS(t-6), PS(t-2), PS(t-8), Qh(t-2), Qh(t-4)$
C2	$PS(t-4), PS(t-6), PS(t-2), PS(t-8), Qh(t-2)$
C3	$PS(t-4), PS(t-6), PS(t-2), PS(t-8)$
C4	$PS(t-4), PS(t-6), PS(t-2)$
C5	$PS(t-4), PS(t-6)$
C6	$PS(t-4)$
C7	$PS(t-2), Qh(t-2)$
C8	$PS(t-4), Qh(t-4)$
C9	$PS(t-6), TLML(t-6)$
C10	$PS(t-8), TLML(t-8)$

4.8

4.9

Figure 7 summarizes the prediction performance of the GBR algorithm for future forecasting of SDMC using the ten input variable combinations in Table 4, evaluated based on the PCC, NSE, N-RMSE and N-MAE indices.

4.10

4.11

4.12

[Figure 7]

4.13

From Figure 7, the prediction performance decreases when moving from lead time ‘t-2’ to ‘t-8’ for the input variables, as also observed in previous studies (Sharafati, Haji Seyed Asadollah, *et al.*, 2020). The C1 input variable combination, including the highest number of input variables (six), is associated with the best prediction performance ($PCC = 0.698$, $NRMSE = 0.733$). However, the C6 input variable combination, including only one variable ($PS(t-4)$) shows only an 8% reduction in accuracy ($PCC = 0.640$, $NRMSE = 0.781$), which is considered an insignificant performance reduction. C6 is therefore the optimal input variable combination, because it only requires the knowledge of a single variable (PS). The preferable use of lead time ‘t-4’ is also confirmed by the high performance of the C8 input variable combination ($PCC = 0.665$, $NRMSE = 0.759$).

4.14

4.15

4.16

4.17

4.18

4.19

4.20

4.21

4.22

4.23

4.24

4. Discussion and Conclusion

Being able to forecast Particulate Matter (PM) concentrations is essential, due to its effects on human life, economy and environment. This study aimed to simulate Surface Mass Dust Concentration using the Modern-Era Retrospective analysis for Research and Applications Version 2 (MERRA-2) aerosol satellite diagnosis. To do this, monthly SDMC data for the period 2009 to 2020 were downloaded from the MERRA-2 database.

The monthly evaluation of air dust distribution in Iran showed that the southeastern regions are characterized by higher dust concentrations, negatively affecting air quality. The province of Khuzestan is the most impacted, as also confirmed by other investigations (Sabetghadam *et al.*, 2018, Yousefi *et al.*, 2020), and was therefore selected as the case study. The year 2018 was specifically considered, because it was characterized by high SDMC and has complete recorded satellite observations. SDMC hourly data were obtained from MERRA-2 for 2018 for the Khuzestan province.

Pearson's Correlation Coefficient (PCC) computation to evaluate the correlation between MERRA-2 SDMC and 11 meteorological, hydrological and geological parameters from Global Precipitation Measurement (GPM), Global Land Data Assimilation System (GLDAS) and another MERRA-2 database showed the wind-related variables - surface wind speed (SPEEDLML), surface aerodynamic conductivity (ACOND) and surface pressure (PS) - to be the most correlated with SDMC. Combinations of these three parameters were evaluated for real-time prediction of SDMC using the Gradient Boosting Regression (GBR) algorithm at 18 Research Points. The best prediction performance was obtained at Research Point 3 ($CC = 0.815$ and $N - RMSE = 0.605$), which was considered for the further forecasting analysis. The input variables are solely based on remote sensing (one of

451 the elements of novelty of this study), therefore the related errors and uncertainties
452 are expected to affect the SDMC prediction performance, compared to models using
453 ground-based measurements. However, while for some Research Points the GBR
454 algorithm produced predictions with PCC value in the moderate acceptance range
455 ($0.3 \leq PCC \leq 0.7$) based on (Ratner, 2009), there are several Research Points where
456 the prediction performance was strong ($0.7 \leq PCC \leq 1.0$). This outcome confirms
457 the prediction potential of ensemble algorithms when using data affected by errors
458 and modelling processes characterized by non-linearity. Comparing this study's
459 results with those by (Nabavi *et al.*, 2018) shows that the GBR algorithm
460 outperforms both SVM ($PCC = 0.81$) and MARS ($PCC = 0.80$). The GBR also
461 appears to have better accuracy than ANN ($PCC = 0.62$), ANFIS ($PCC = 0.70$) and
462 GRNN ($PCC = 0.71$) (Mirzaei *et al.*, 2019). From a spatial point of view, our results
463 highlighted a better prediction performance for the southern low lands of the
464 Khuzestan province, when compared with the higher and mountainous lands in the
465 north. Predictions were also better in marshlands compared with rocky soils.

466 While the wind-related input variables govern the real time ('t') prediction of
467 SDMC, the heat-related variables are also important for the future forecasting of
468 SDMC (lead times of 't-2' to 't-8'). Considering Research Point 3 for the analysis,
469 the PS variable allows for the better forecasting, closely followed by surface upward
470 sensible heat flux (Qh) and surface air temperature (TLML). It is worth mentioning
471 that, as the lead time goes from 't-2' to 't-8', the influence of TLML on SDMC
472 forecasting becomes stronger than that of Qh. The evaluation of input variable
473 combinations for future SDMC forecasting revealed that the use of a single input
474 variable, *PS* (t-4), with $PCC = 0.640$ and $N - RMSE = 0.781$, is the optimal
475 approach (most cost-efficient and applicable). To the authors' knowledge, there are
476 no previous studies specifically aimed at forecasting future air pollutant

477 concentrations. The future forecasting performance obtained here with the GBR
478 algorithm is comparable with previously presented real-time predictions.

479

480 **Declaration of interests**

481 The authors declare that they have no known competing financial interests or
482 personal relationships that could have appeared to influence the work reported in
483 this paper.

484

485

486 **References**

487 Al-Othman, A., Tawalbeh, M., Martis, R., Dhou, S., Orhan, M., Qasim, M. &
488 Olabi, A. G. (2022) Artificial intelligence and numerical models in hybrid
489 renewable energy systems with fuel cells: Advances and prospects. *Energy*
490 *Conversion and Management* **253**, 115154. Elsevier.

491 Amaral, S. S., Carvalho, J. A. De, Costa, M. A. M. & Pinheiro, C. (2015) An
492 overview of particulate matter measurement instruments. *Atmosphere* **6**(9),
493 1327–1345. Multidisciplinary Digital Publishing Institute.

494 Anderson, J. O., Thundiyil, J. G. & Stolbach, A. (2012) Clearing the air: a review
495 of the effects of particulate matter air pollution on human health. *Journal of*
496 *medical toxicology* **8**(2), 166–175. Springer.

497 Arkian, F. & Nicholson, S. E. (2018) Long-term variations of aerosol optical depth
498 and aerosol radiative forcing over Iran based on satellite and AERONET data.
499 *Environmental monitoring and assessment* **190**(1), 1–15. Springer.

- 000 Asadollah, S. B. H. S., Khan, N., Sharafati, A., Shahid, S., Chung, E.-S. & Wang,
001 X.-J. (2021) Prediction of heat waves using meteorological variables in diverse
002 regions of Iran with advanced machine learning models. *Stochastic*
003 *Environmental Research and Risk Assessment* 1–16. Springer.
- 004 Brook, R. D. & Rajagopalan, S. (2009) Particulate matter, air pollution, and blood
005 pressure. *Journal of the American Society of Hypertension* 3(5), 332–350.
006 Elsevier.
- 007 Campbell, J. B. & Wynne, R. H. (2011) *Introduction to remote sensing*. Guilford
008 Press.
- 009 Chen, T.-M., Kuschner, W. G., Gokhale, J. & Shofer, S. (2007) Outdoor air
010 pollution: nitrogen dioxide, sulfur dioxide, and carbon monoxide health
011 effects. *The American journal of the medical sciences* 333(4), 249–256.
012 Elsevier.
- 013 Chok, N. S. (2010) Pearson’s versus Spearman’s and Kendall’s correlation
014 coefficients for continuous data. University of Pittsburgh.
- 015 Chu, Y., Liu, Y., Li, X., Liu, Z., Lu, H., Lu, Y., Mao, Z., et al. (2016) A review on
016 predicting ground PM_{2.5} concentration using satellite aerosol optical depth.
017 *Atmosphere* 7(10), 129. Multidisciplinary Digital Publishing Institute.
- 018 Chudnovsky, A. A., Koutrakis, P., Kloog, I., Melly, S., Nordio, F., Lyapustin, A.,
019 Wang, Y., et al. (2014) Fine particulate matter predictions using high
020 resolution Aerosol Optical Depth (AOD) retrievals. *Atmospheric Environment*
021 89, 189–198. Elsevier.
- 022 Chudnovsky, A., Tang, C., Lyapustin, A., Wang, Y., Schwartz, J. & Koutrakis, P.
023 (2013) A critical assessment of high-resolution aerosol optical depth retrievals

- 024 for fine particulate matter predictions. *Atmospheric Chemistry and Physics*
025 **13**(21), 10907–10917. Copernicus GmbH.
- 026 Council, N. R. (1992) *Rethinking the ozone problem in urban and regional air*
027 *pollution*. National Academies Press.
- 028 Dadashi-Roudbari, A. & Ahmadi, M. (2020) Evaluating temporal and spatial
029 variability and trend of aerosol optical depth (550 nm) over Iran using data
030 from MODIS on board the Terra and Aqua satellites. *Arabian Journal of*
031 *Geosciences* **13**(6), 1–23. Springer.
- 032 Daniali, M. & Karimi, N. (2019) Spatiotemporal analysis of dust patterns over
033 Mesopotamia and their impact on Khuzestan province, Iran. *Natural Hazards*
034 **97**(1), 259–281. Springer.
- 035 Davis, S. M. & Swain, P. H. (1978) *Remote sensing: the quantitative approach*.
036 McGraw-Hill International Book Company New York.
- 037 Di, Q., Kloog, I., Koutrakis, P., Lyapustin, A., Wang, Y. & Schwartz, J. (2016)
038 Assessing PM_{2.5} exposures with high spatiotemporal resolution across the
039 continental United States. *Environmental science & technology* **50**(9), 4712–
040 4721. ACS Publications.
- 041 Diao, M., Holloway, T., Choi, S., O’Neill, S. M., Al-Hamdan, M. Z., Donkelaar,
042 A. Van, Martin, R. V, et al. (2019) Methods, availability, and applications of
043 PM_{2.5} exposure estimates derived from ground measurements, satellite, and
044 atmospheric models. *Journal of the Air & Waste Management Association*
045 **69**(12), 1391–1414. Taylor & Francis.
- 046 Donkelaar, A. Van, Martin, R. V, Brauer, M., Kahn, R., Levy, R., Verduzco, C. &
047 Villeneuve, P. J. (2010) Global estimates of ambient fine particulate matter

- 058 concentrations from satellite-based aerosol optical depth: development and
059 application. *Environmental health perspectives* **118**(6), 847–855. National
060 Institute of Environmental Health Sciences.
- 061 Duan, Y., Liao, Y., Li, H., Yan, S., Zhao, Z., Yu, S., Fu, Y., et al. (2019) Effect of
062 changes in season and temperature on cardiovascular mortality associated with
063 nitrogen dioxide air pollution in Shenzhen, China. *Science of The Total
064 Environment* **697**, 134051. Elsevier.
- 065 Friedman, J. H. (2001) Greedy function approximation: a gradient boosting
066 machine. *Annals of statistics* 1189–1232. JSTOR.
- 067 Gelaro, R., McCarty, W., Suárez, M. J., Todling, R., Molod, A., Takacs, L.,
068 Randles, C. A., et al. (2017) The modern-era retrospective analysis for
069 research and applications, version 2 (MERRA-2). *Journal of climate* **30**(14),
070 5419–5454.
- 071 Ghozat, A., Sharafati, A. & Hosseini, S. A. (2022) Satellite-based monitoring of
072 meteorological drought over different regions of Iran: application of the
073 CHIRPS precipitation product. *Environmental Science and Pollution Research*
074 1–18. Springer.
- 075 Gomis, D., Ruiz, S., Sotillo, M. G., Álvarez-Fanjul, E. & Terradas, J. (2008) Low
076 frequency Mediterranean sea level variability: the contribution of atmospheric
077 pressure and wind. *Global and Planetary Change* **63**(2–3), 215–229. Elsevier.
- 078 Gueymard, C. A. & Yang, D. (2020) Worldwide validation of CAMS and
079 MERRA-2 reanalysis aerosol optical depth products using 15 years of
080 AERONET observations. *Atmospheric Environment* **225**, 117216. Elsevier.
- 081 Guo, H., Xu, M. & Hu, Q. (2011) Changes in near-surface wind speed in China:

- 052 1969–2005. *International Journal of Climatology* **31**(3), 349–358. Wiley
053 Online Library.
- 054 Hamanaka, R. B. & Mutlu, G. M. (2018) Particulate matter air pollution: effects on
055 the cardiovascular system. *Frontiers in endocrinology* **9**, 680. Frontiers.
- 056 Hauke, J. & Kossowski, T. (2011) Comparison of values of Pearson's and
057 Spearman's correlation coefficient on the same sets of data. Wydział Nauk
058 Geograficznych i Geologicznych Uniwersytetu im. Adama Mickiewicza.
- 059 Jodar-Abellan, A., Valdes-Abellan, J., Pla, C. & Gomariz-Castillo, F. (2019)
060 Impact of land use changes on flash flood prediction using a sub-daily SWAT
061 model in five Mediterranean ungauged watersheds (SE Spain). *Science of the*
062 *Total Environment* **657**, 1578–1591. Elsevier.
- 063 Johnson, N. E., Bonczak, B. & Kontokosta, C. E. (2018) Using a gradient boosting
064 model to improve the performance of low-cost aerosol monitors in a dense,
065 heterogeneous urban environment. *Atmospheric environment* **184**, 9–16.
066 Elsevier.
- 067 Just, A. C., Wright, R. O., Schwartz, J., Coull, B. A., Baccarelli, A. A., Tellez-
068 Rojo, M. M., Moody, E., et al. (2015) Using high-resolution satellite aerosol
069 optical depth to estimate daily PM_{2.5} geographical distribution in Mexico
070 City. *Environmental science & technology* **49**(14), 8576–8584. ACS
071 Publications.
- 072 Karandish, F. & Šimůnek, J. (2016) A comparison of numerical and machine-
073 learning modeling of soil water content with limited input data. *Journal of*
074 *Hydrology* **543**, 892–909. Elsevier.
- 075 Khanzode, K. C. A. & Sarode, R. D. (2020) Advantages and Disadvantages of

- 096 Artificial Intelligence and Machine Learning: A Literature Review.
097 *International Journal of Library & Information Science (IJLIS)* **9**(1), 3.
- 098 Kianian, B., Liu, Y. & Chang, H. H. (2021) Imputing satellite-derived aerosol
099 optical depth using a multi-resolution spatial model and random forest for
100 PM_{2.5} prediction. *Remote Sensing* **13**(1), 126. Multidisciplinary Digital
101 Publishing Institute.
- 102 Kwasny, F., Madl, P. & Hofmann, W. (2010) Correlation of air quality data to
103 ultrafine particles (UFP) concentration and size distribution in ambient air.
104 *Atmosphere* **1**(1), 3–14. Molecular Diversity Preservation International.
- 105 Lary, D. J., Faruque, F. S., Malakar, N., Moore, A., Roscoe, B., Adams, Z. L. &
106 Eggelston, Y. (2014) Estimating the global abundance of ground level
107 presence of particulate matter (PM_{2.5}). *Geospatial health* **8**(3), S611–S630.
- 108 Lee, H. J., Coull, B. A., Bell, M. L. & Koutrakis, P. (2012) Use of satellite-based
109 aerosol optical depth and spatial clustering to predict ambient PM_{2.5}
110 concentrations. *Environmental research* **118**, 8–15. Elsevier.
- 111 Lee, M., Kloog, I., Chudnovsky, A., Lyapustin, A., Wang, Y., Melly, S., Coull, B.,
112 et al. (2016) Spatiotemporal prediction of fine particulate matter using high-
113 resolution satellite images in the Southeastern US 2003–2011. *Journal of*
114 *exposure science & environmental epidemiology* **26**(4), 377–384. Nature
115 Publishing Group.
- 116 Liu, D. & Li, L. (2015) Application study of comprehensive forecasting model
117 based on entropy weighting method on trend of PM_{2.5} concentration in
118 Guangzhou, China. *International journal of environmental research and*
119 *public health* **12**(6), 7085–7099. Multidisciplinary Digital Publishing Institute.

- 720 Liu, S., Lu, L., Mao, D. & Jia, L. (2007) Evaluating parameterizations of
721 aerodynamic resistance to heat transfer using field measurements. *Hydrology
722 and earth system sciences* **11**(2), 769–783. Copernicus GmbH.
- 723 Mallick, K., Wandera, L., Bhattarai, N., Hostache, R., Kleniewska, M. &
724 Chormanski, J. (2018) A critical evaluation on the role of aerodynamic and
725 canopy–surface conductance parameterization in SEB and SVAT models for
726 simulating evapotranspiration: A case study in the upper bieberza national park
727 wetland in poland. *Water* **10**(12), 1753. Multidisciplinary Digital Publishing
728 Institute.
- 729 Mehdizadeh, S., Behmanesh, J. & Khalili, K. (2017) Evaluating the performance
730 of artificial intelligence methods for estimation of monthly mean soil
731 temperature without using meteorological data. *Environmental Earth Sciences*
732 **76**(8), 1–16. Springer.
- 733 MERRA, G. (2AD) tavgU_2d_lnd_Nx: 2d, diurnal, time-averaged, single-level,
734 assimilation, land surface diagnostics V5. 12.4. *EarthData GES DISC NASA*.
- 735 Mirakbari, M. & Ebrahimi Khusfi, Z. (2020) Investigation of spatial and temporal
736 changes in atmospheric aerosol using aerosol optical depth in Southeastern
737 Iran. *Journal of RS and GIS for Natural Resources* **11**(3), 87–105. Bushehr
738 Branch, Islamic Azad University.
- 739 Mirzaei, M., Amanollahi, J. & Tzani, C. G. (2019) Evaluation of linear, nonlinear,
740 and hybrid models for predicting PM 2.5 based on a GTWR model and
741 MODIS AOD data. *Air Quality, Atmosphere & Health* **12**(10), 1215–1224.
742 Springer.
- 743 Moriasi, D. N., Arnold, J. G., Liew, M. W. Van, Bingner, R. L., Harmel, R. D. &
744 Veith, T. L. (2007) Model evaluation guidelines for systematic quantification

- 745 of accuracy in watershed simulations. *Transactions of the ASABE* **50**(3), 885–
746 900. American society of agricultural and biological engineers.
- 747 Müller, A. C. & Guido, S. (2016) *Introduction to machine learning with Python: a*
748 *guide for data scientists*. ‘ O’Reilly Media, Inc.’
- 749 Nabavi, S. O., Haimberger, L., Abbasi, R. & Samimi, C. (2018) Prediction of
750 aerosol optical depth in West Asia using deterministic models and machine
751 learning algorithms. *Aeolian research* **35**, 69–84. Elsevier.
- 752 Nakata, M., Sano, I., Mukai, S. & Holben, B. N. (2013) Spatial and temporal
753 variations of atmospheric aerosol in Osaka. *Atmosphere* **4**(2), 157–168.
754 Multidisciplinary Digital Publishing Institute.
- 755 Nguyen, D. L., Kim, J. Y., Ghim, Y. S. & Shim, S.-G. (2015) Influence of regional
756 biomass burning on the highly elevated organic carbon concentrations
757 observed at Gosan, South Korea during a strong Asian dust period.
758 *Environmental Science and Pollution Research* **22**(5), 3594–3605. Springer.
- 759 Nie, P., Roccotelli, M., Fanti, M. P., Ming, Z. & Li, Z. (2021) Prediction of home
760 energy consumption based on gradient boosting regression tree. *Energy*
761 *Reports* **7**, 1246–1255. Elsevier.
- 762 Nourani, V., Baghanam, A. H., Adamowski, J. & Kisi, O. (2014) Applications of
763 hybrid wavelet–artificial intelligence models in hydrology: a review. *Journal*
764 *of Hydrology* **514**, 358–377. Elsevier.
- 765 Pardo, M. Á., Riquelme, A. J., Jodar-Abellan, A. & Melgarejo, J. (2020) Water and
766 energy demand management in pressurized irrigation networks. *Water* **12**(7),
767 1878. Multidisciplinary Digital Publishing Institute.
- 768 Planning, U. S. E. P. A. O. of A. Q. (1996) *Review of the national ambient air*

- 679 *quality standards for particulate matter: Policy assessment of scientific and*
670 *technical information.* DIANE Publishing.
- 671 Raaschou-Nielsen, O., Beelen, R., Wang, M., Hoek, G., Andersen, Z. J.,
672 Hoffmann, B., Stafoggia, M., et al. (2016) Particulate matter air pollution
673 components and risk for lung cancer. *Environment international* **87**, 66–73.
674 Elsevier.
- 675 Randles, C. A., Silva, A. M. Da, Buchard, V., Colarco, P. R., Darmenov, A.,
676 Govindaraju, R., Smirnov, A., et al. (2017) The MERRA-2 aerosol reanalysis,
677 1980 onward. Part I: System description and data assimilation evaluation.
678 *Journal of climate* **30**(17), 6823–6850.
- 679 Ratner, B. (2009) The correlation coefficient: Its values range between + 1/– 1, or
680 do they? *Journal of targeting, measurement and analysis for marketing* **17**(2),
681 139–142. Springer.
- 682 Rebekić, A., Lončarić, Z., Petrović, S. & Marić, S. (2015) Pearson's or Spearman's
683 correlation coefficient-which one to use? *Poljoprivreda* **21**(2), 47–54. Fakultet
684 agrobiotehničkih znanosti Osijek i Poljoprivredni institut Osijek.
- 685 Reichle, R. H., Draper, C. S., Liu, Q., Giroto, M., Mahanama, S. P. P., Koster, R.
686 D. & Lannoy, G. J. M. De. (2017) Assessment of MERRA-2 land surface
687 hydrology estimates. *Journal of Climate* **30**(8), 2937–2960.
- 688 Rezaei, M., Farajzadeh, M., Mielonen, T. & Ghavidel, Y. (2019) Analysis of
689 spatio-temporal dust aerosol frequency over Iran based on satellite data.
690 *Atmospheric Pollution Research* **10**(2), 508–519. Elsevier.
- 691 Rodell, M., Houser, P. R., Jambor, U. E. A., Gottschalck, J., Mitchell, K., Meng,
692 C.-J. J., Arsenault, K., et al. (2004) The global land data assimilation system.

- ٦٩٣ *Bulletin of the American Meteorological Society* **85**(3), 381–394. American
٦٩٤ Meteorological Society. doi:10.1175/BAMS-85-3-381
- ٦٩٥ Sabetghadam, S., Khoshsima, M. & Alizadeh-Choobari, O. (2018) Spatial and
٦٩٦ temporal variations of satellite-based aerosol optical depth over Iran in
٦٩٧ Southwest Asia: Identification of a regional aerosol hot spot. *Atmospheric
٦٩٨ Pollution Research* **9**(5), 849–856. Elsevier.
- ٦٩٩ Salami, H., Khorami, S., Yazdani, S. & Saleh, I. (2021) Economic Evaluation of
٧٠٠ the Damages of Dust Bowl on Crop Yield by Choice Experiment Method in
٧٠١ Khuzestan Province of Iran. *International Journal of Agricultural
٧٠٢ Management and Development* **11**(3). Islamic Azad University, Rasht Branch.
- ٧٠٣ Sharafati, A., Asadollah, S. B. H. S. & Hosseinzadeh, M. (2020) The potential of
٧٠٤ new ensemble machine learning models for effluent quality parameters
٧٠٥ prediction and related uncertainty. *Process Safety and Environmental
٧٠٦ Protection*. Elsevier.
- ٧٠٧ Sharafati, A., Asadollah, S. B. H. S. & Neshat, A. (2020) A new artificial
٧٠٨ intelligence strategy for predicting the groundwater level over the Rafsanjan
٧٠٩ aquifer in Iran. *Journal of Hydrology* **591**, 125468. Elsevier.
- ٧١٠ Sharafati, A., Haji Seyed Asadollah, S. B., Motta, D. & Yaseen, Z. M. (2020)
٧١١ Application of newly developed ensemble machine learning models for daily
٧١٢ suspended sediment load prediction and related uncertainty analysis.
٧١٣ *Hydrological Sciences Journal*. Taylor & Francis.
- ٧١٤ Shiru, M. S., Shahid, S., Chae, S.-T. & Chung, E.-S. (2022) Replicability of
٧١٥ Annual and Seasonal Precipitation by CMIP5 and CMIP6 GCMs over East
٧١٦ Asia. *KSCE Journal of Civil Engineering* 1–12. Springer.

- 717 Srivastava, C., Singh, S. & Singh, A. P. (2018) Estimation of air pollution in Delhi
718 using machine learning techniques. *2018 International Conference on*
719 *Computing, Power and Communication Technologies (GUCON)*, 304–309.
720 IEEE.
- 721 Sun, E., Xu, X., Che, H., Tang, Z., Gui, K., An, L., Lu, C., et al. (2019) Variation
722 in MERRA-2 aerosol optical depth and absorption aerosol optical depth over
723 China from 1980 to 2017. *Journal of Atmospheric and Solar-Terrestrial*
724 *Physics* **186**, 8–19. Elsevier.
- 725 Sunyer, J., Ballester, F., Tertre, A. Le, Atkinson, R., Ayres, J. G., Forastiere, F.,
726 Forsberg, B., et al. (2003) The association of daily sulfur dioxide air pollution
727 levels with hospital admissions for cardiovascular diseases in Europe (The
728 Aphea-II study). *European heart journal* **24**(8), 752–760. Oxford University
729 Press.
- 730 Ukhov, A., Mostamandi, S., Silva, A. Da, Flemming, J., Alshehri, Y., Shevchenko,
731 I. & Stenchikov, G. (2020) Assessment of natural and anthropogenic aerosol
732 air pollution in the Middle East using MERRA-2, CAMS data assimilation
733 products, and high-resolution WRF-Chem model simulations. *Atmospheric*
734 *Chemistry and Physics* **20**(15), 9281–9310. Copernicus GmbH.
- 735 USEPA, O. A. R. (2019) Health and environmental effects of particulate matter
736 (PM). Retrieved.
- 737 Veselovskii, I., Goloub, P., Podvin, T., Tanre, D., Silva, A. Da, Colarco, P.,
738 Castellanos, P., et al. (2018) Characterization of smoke and dust episode over
739 West Africa: comparison of MERRA-2 modeling with multiwavelength Mie–
740 Raman lidar observations. *Atmospheric measurement techniques* **11**(2), 949–
741 969. Copernicus GmbH.

- 742 Wang, F., Wang, Y., Zhang, K., Hu, M., Weng, Q. & Zhang, H. (2021) Spatial
743 heterogeneity modeling of water quality based on random forest regression and
744 model interpretation. *Environmental Research* **202**, 111660. Elsevier.
- 745 Wang, W.-C., Chau, K.-W., Cheng, C.-T. & Qiu, L. (2009) A comparison of
746 performance of several artificial intelligence methods for forecasting monthly
747 discharge time series. *Journal of hydrology* **374**(3–4), 294–306. Elsevier.
- 748 WHO. (2014) 7 million premature deaths annually linked to air pollution.
749 Retrieved March 25, 2014, from [https://www.who.int/news/item/25-03-2014-](https://www.who.int/news/item/25-03-2014-7-million-premature-deaths-annually-linked-to-air-pollution)
750 [7-million-premature-deaths-annually-linked-to-air-pollution](https://www.who.int/news/item/25-03-2014-7-million-premature-deaths-annually-linked-to-air-pollution)
- 751 Wuest, T., Weimer, D., Irgens, C. & Thoben, K.-D. (2016) Machine learning in
752 manufacturing: advantages, challenges, and applications. *Production &*
753 *Manufacturing Research* **4**(1), 23–45. Taylor & Francis.
- 754 Xu, X., Wu, H., Yang, X. & Xie, L. (2020) Distribution and transport
755 characteristics of dust aerosol over Tibetan Plateau and Taklimakan Desert in
756 China using MERRA-2 and CALIPSO data. *Atmospheric Environment* **237**,
757 117670. Elsevier.
- 758 Yao, W., Che, H., Gui, K., Wang, Y. & Zhang, X. (2020) Can MERRA-2
759 reanalysis data reproduce the three-dimensional evolution characteristics of a
760 typical dust process in East Asia? A case study of the dust event in May 2017.
761 *Remote Sensing* **12**(6), 902. Multidisciplinary Digital Publishing Institute.
- 762 Yousefi, R., Wang, F., Ge, Q. & Shaheen, A. (2020) Long-term aerosol optical
763 depth trend over Iran and identification of dominant aerosol types. *Science of*
764 *The Total Environment* **722**, 137906. Elsevier.
- 765 Zarasvandi, A. (2009) Environmental impacts of dust storms in the Khuzestan

- ۷۶۶ province. *Environmental Protection Agency (EPA) of Khuzestan province,*
۷۶۷ *Internal Report, 375p.*
- ۷۶۸ Zarasvandi, A., Carranza, E. J. M., Moore, F. & Rastmanesh, F. (2011) Spatio-
۷۶۹ temporal occurrences and mineralogical–geochemical characteristics of
۷۷۰ airborne dusts in Khuzestan Province (southwestern Iran). *Journal of*
۷۷۱ *geochemical exploration* **111**(3), 138–151. Elsevier.
- ۷۷۲ Zhang, J., Ma, G., Huang, Y., Aslani, F. & Nener, B. (2019) Modelling uniaxial
۷۷۳ compressive strength of lightweight self-compacting concrete using random
۷۷۴ forest regression. *Construction and Building Materials* **210**, 713–719.
۷۷۵ Elsevier.
- ۷۷۶ Zhang, X., Chu, Y., Wang, Y. & Zhang, K. (2018) Predicting daily PM2. 5
۷۷۷ concentrations in Texas using high-resolution satellite aerosol optical depth.
۷۷۸ *Science of the Total Environment* **631**, 904–911. Elsevier.
- ۷۷۹ Zhang, Y. & Haghani, A. (2015) A gradient boosting method to improve travel
۷۸۰ time prediction. *Transportation Research Part C: Emerging Technologies* **58**,
۷۸۱ 308–324. Elsevier.
- ۷۸۲ Zhu, X., Zhang, P. & Xie, M. (2021) A Joint Long Short-Term Memory and
۷۸۳ AdaBoost regression approach with application to remaining useful life
۷۸۴ estimation. *Measurement* **170**, 108707. Elsevier.

۷۸۵

۷۸۶

۷۸۷

۷۸۸

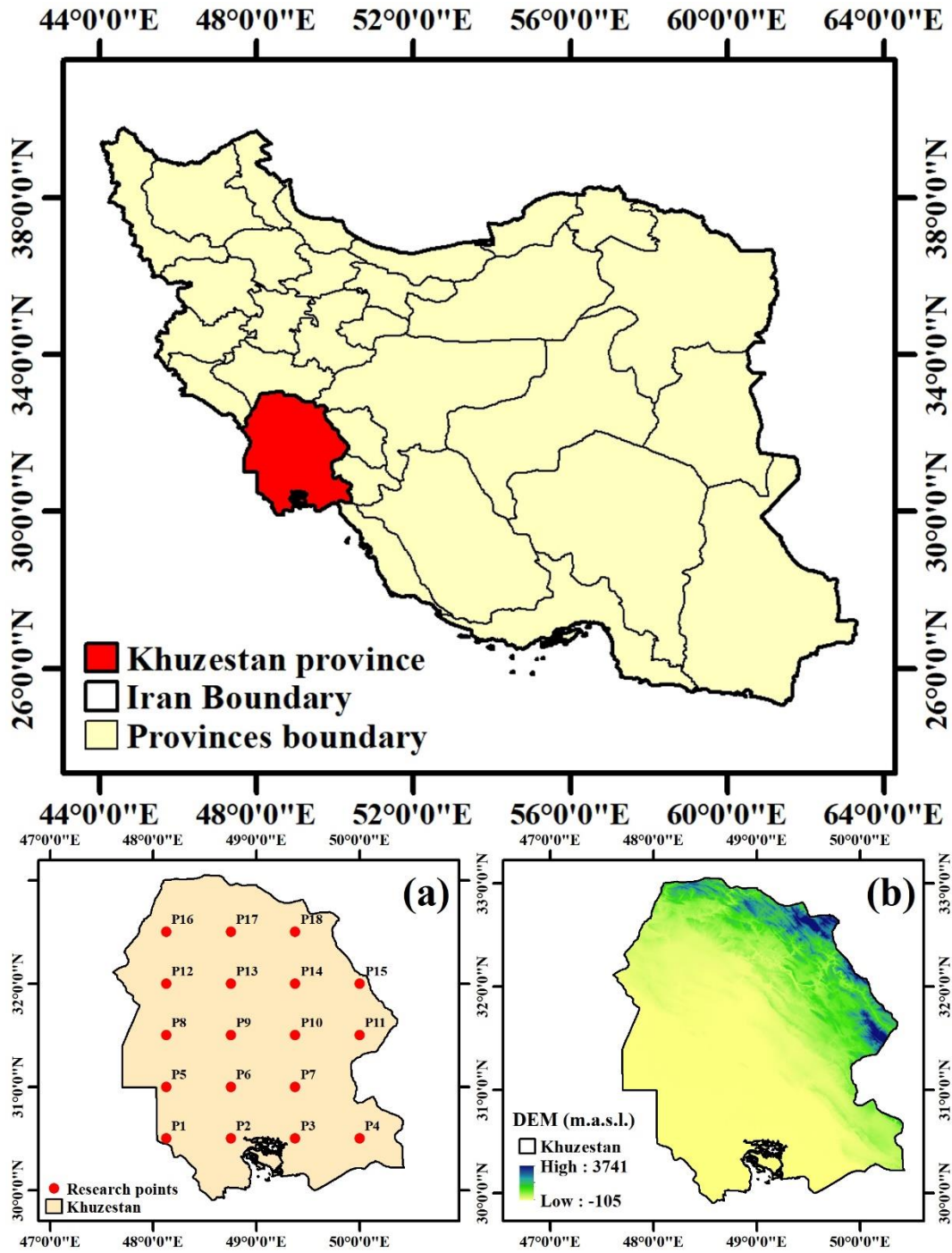


Figure 1: Location of the Khuzestan province in Iran with (a) Research Point distribution and (b) Digital Elevation Model (DEM).

۷۸۹

۷۹۰

۷۹۱

۷۹۲

۷۹۳

۷۹۴

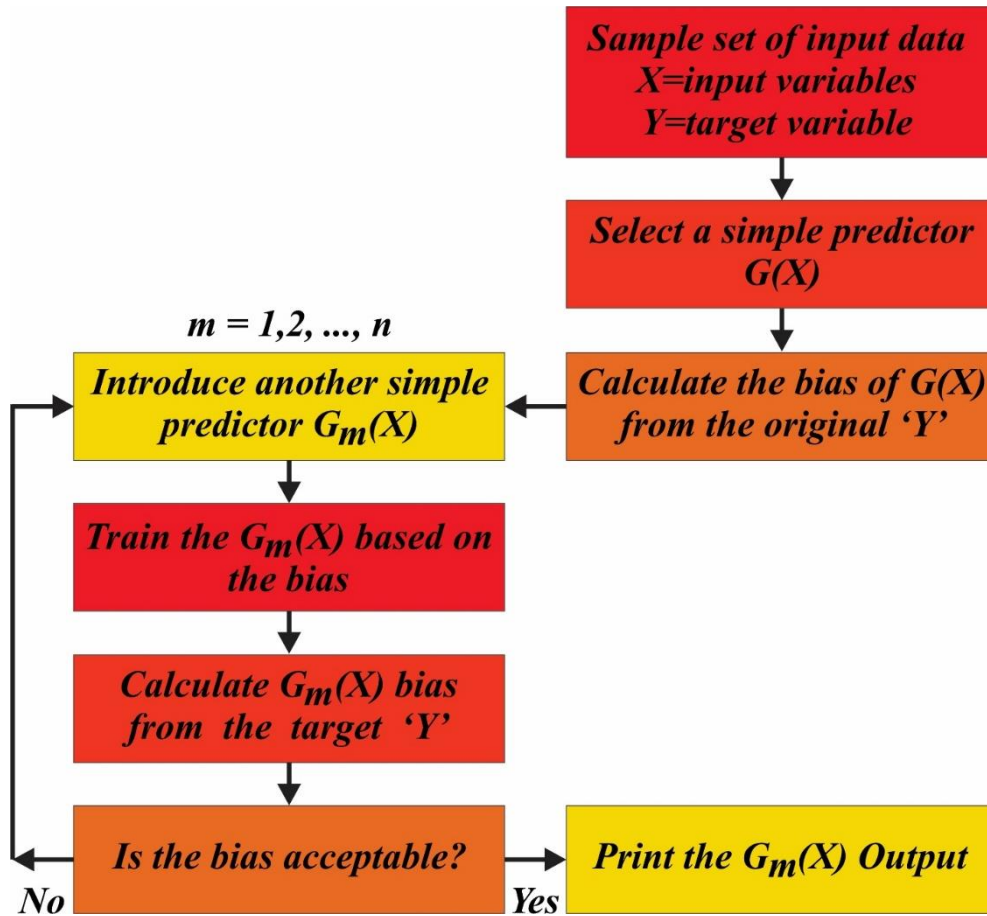
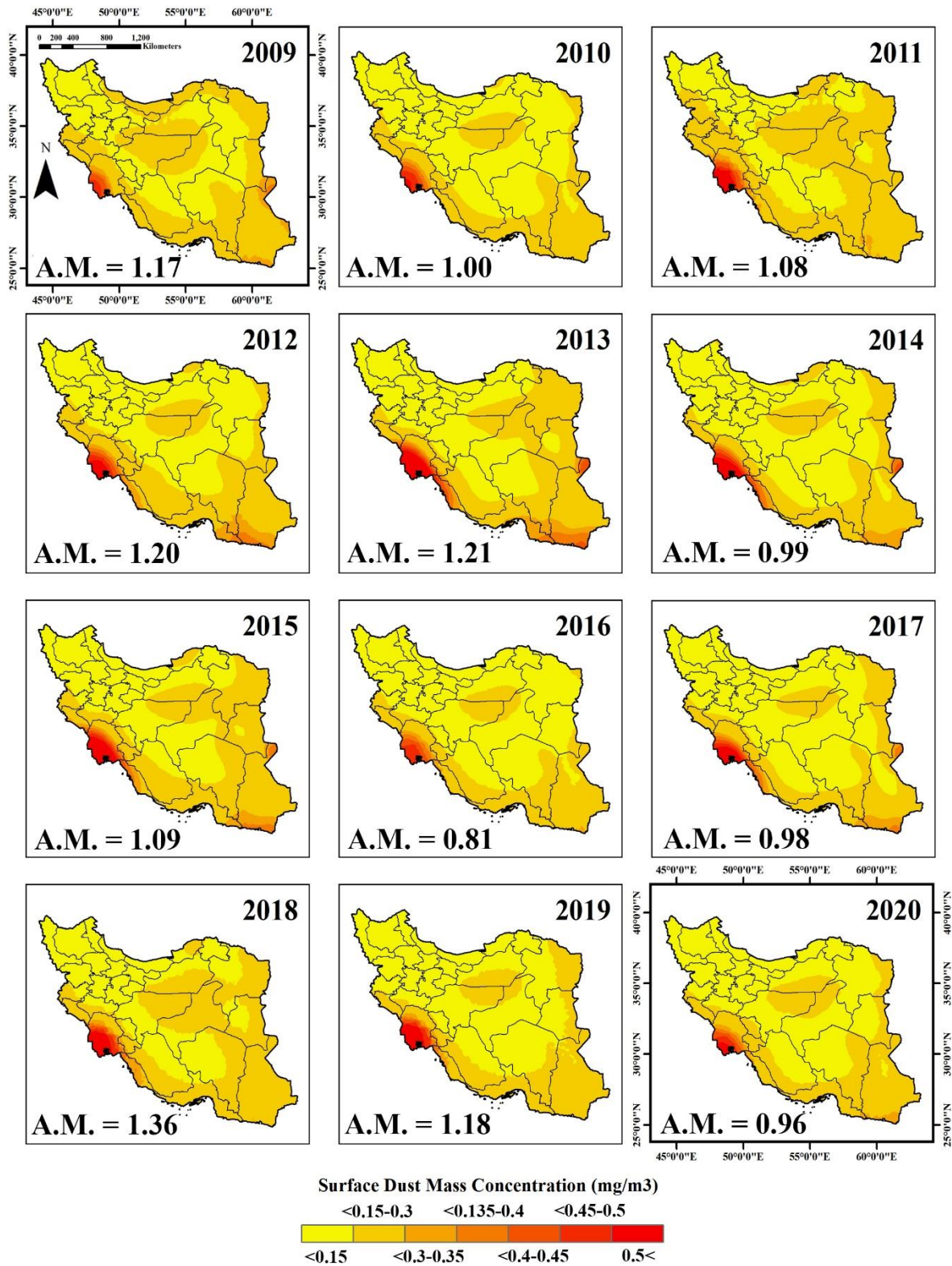


Figure 2: Flowchart of Gradient Boosting Regression.

790
 796
 797
 798
 799
 800
 801
 802
 803



۸.۴

۸.۵

۸.۶

Figure 3: Yearly maximum Surface Dust Mass Concentration (SDMC) between 2009 and 2020 in Iran.

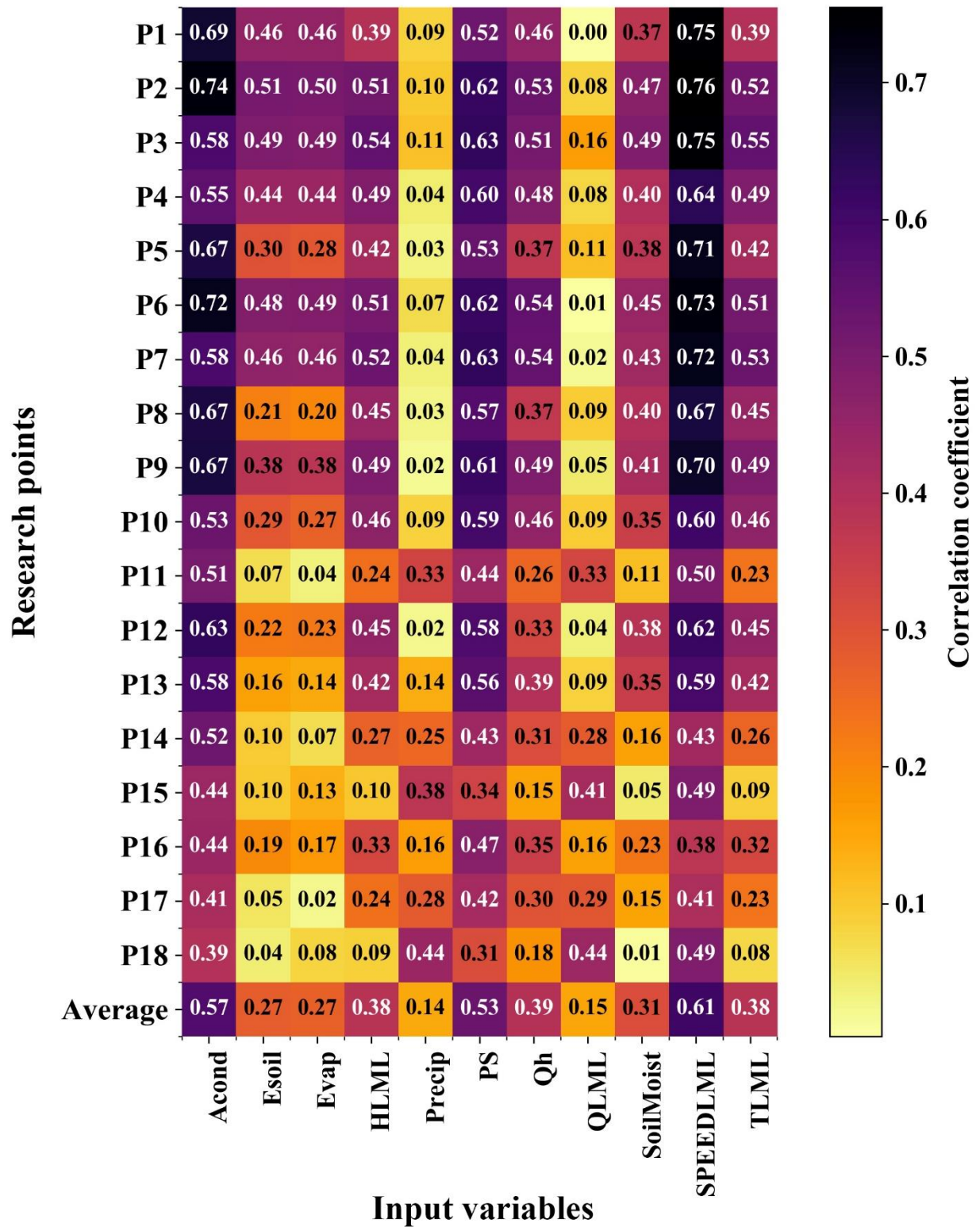
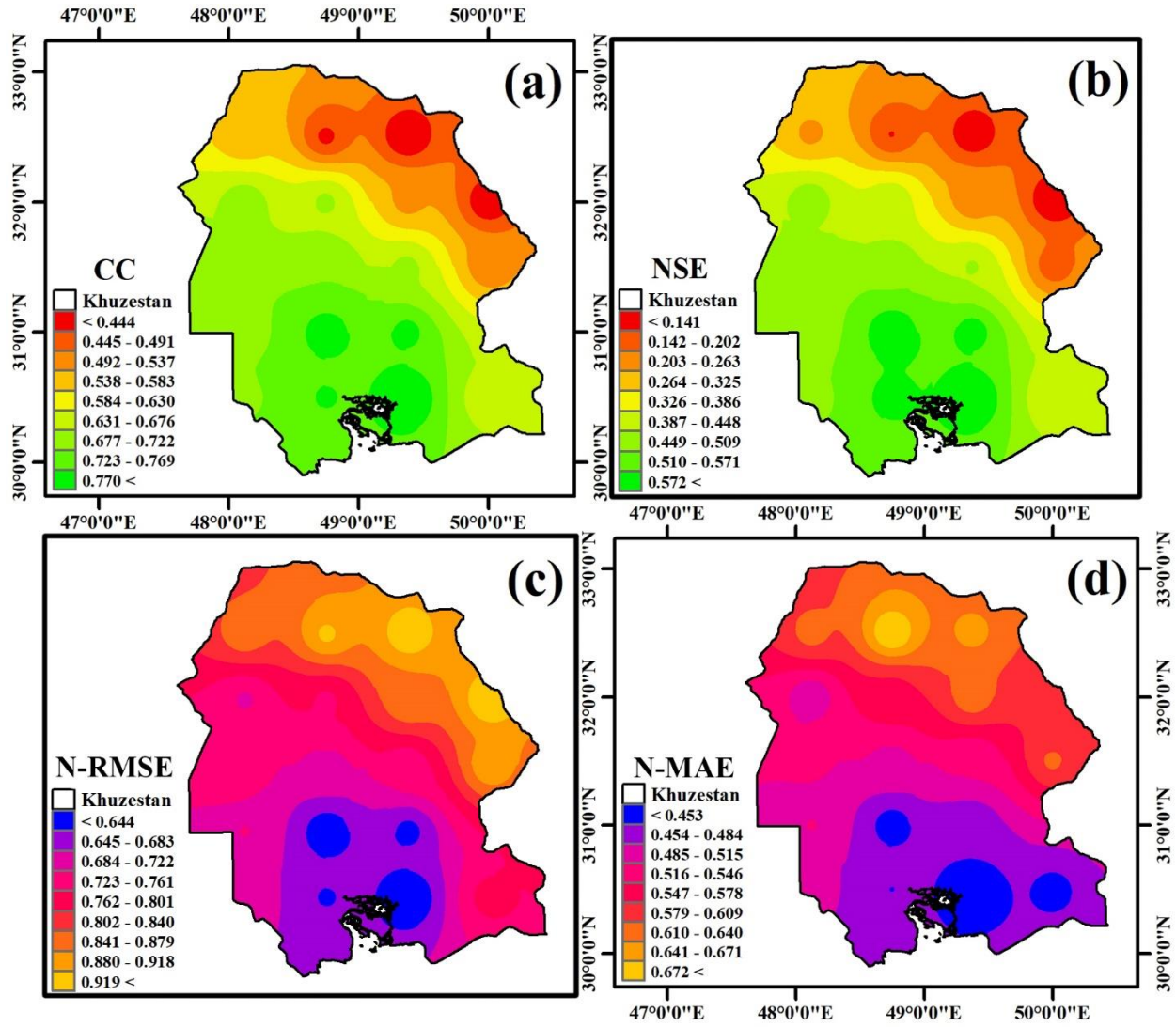


Figure 4: Heat map of Pearson's Correlation Coefficient (PCC) between the 11 input variables considered and daily SDMC for the 18 Research Points.

۸۱۰



۸۱۱

۸۱۲

Figure 5: SDMC real-time prediction performance over the Khuzestan province.

۸۱۳

۸۱۴

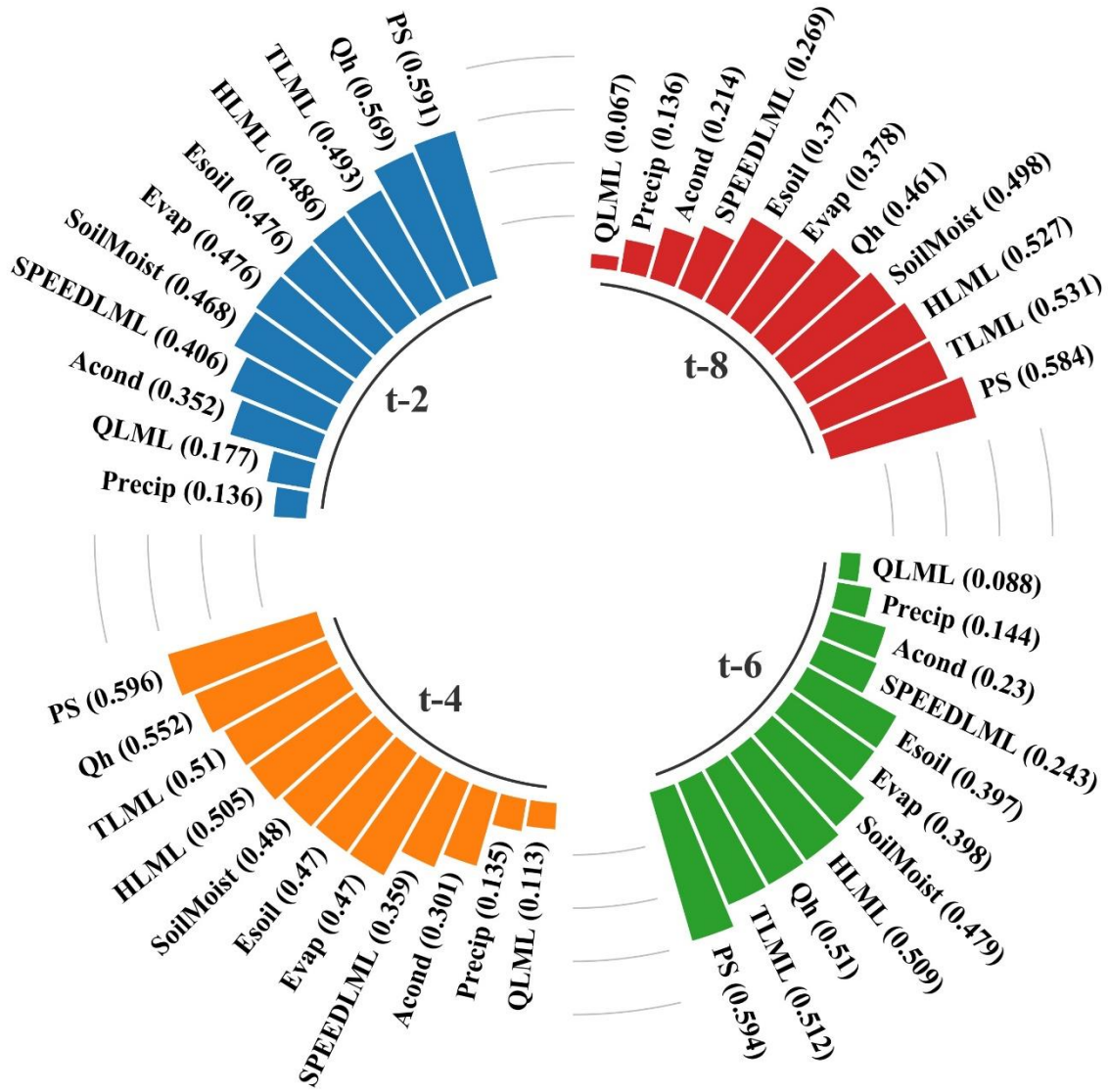
۸۱۵

۸۱۶

۸۱۷

۸۱۸

۸۱۹



۸۲۰
 ۸۲۱
 ۸۲۲

Figure 6: Pearson's Correlation Coefficient for correlation between the 11 input variables considered and SDMC for different lead times.

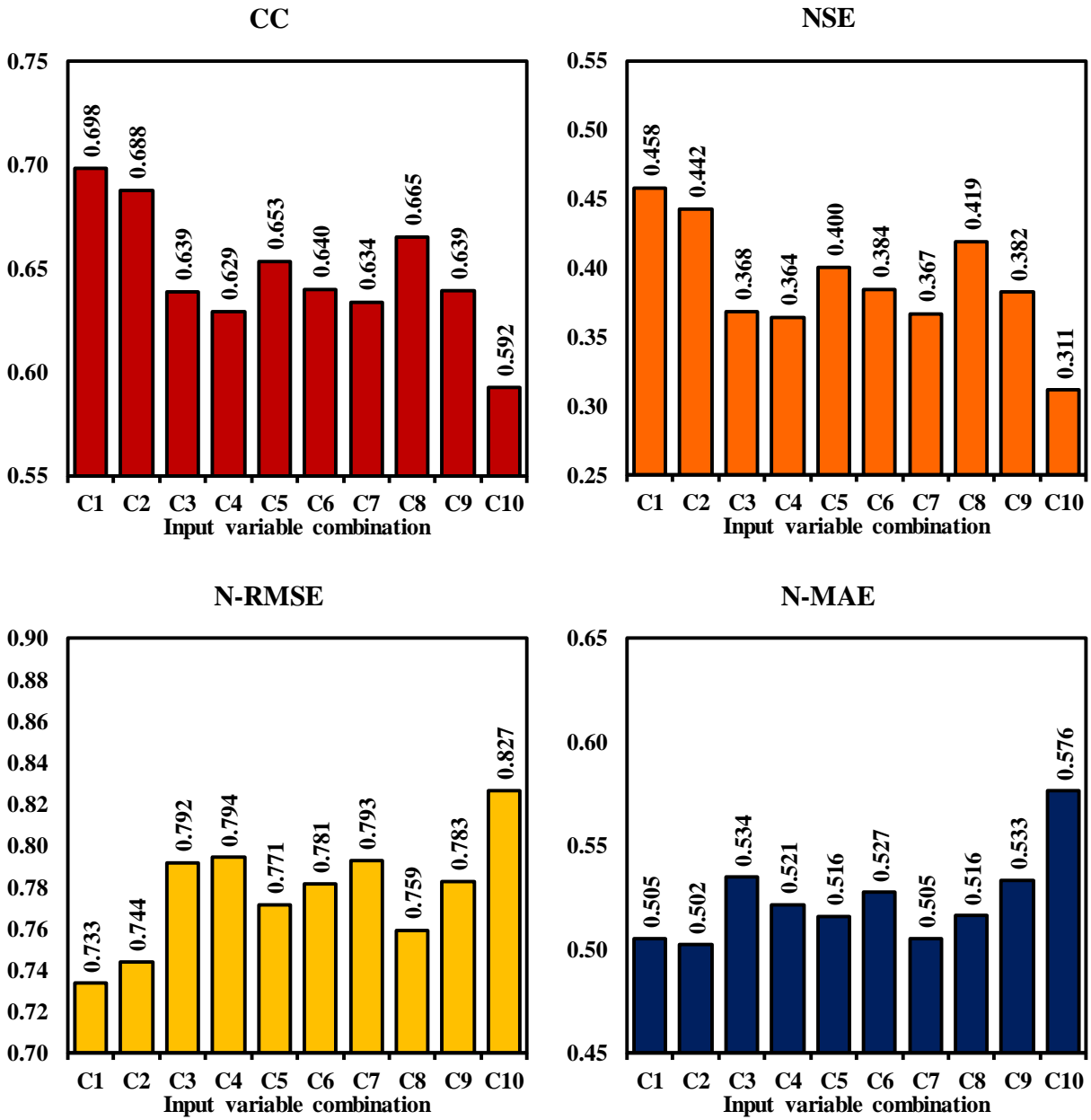


Figure 7: Prediction performance for future forecasting of SDMC for the ten input variable combinations considered.