



# Estimation of daily PM<sub>10</sub> and PM<sub>2.5</sub> concentrations in Italy, 2013–2015, using a spatiotemporal land-use random-forest model

Massimo Stafoggia<sup>a,b,\*</sup>, Tom Bellander<sup>b</sup>, Simone Bucci<sup>a</sup>, Marina Davoli<sup>a</sup>, Kees de Hoogh<sup>c,d</sup>, Francesca de' Donato<sup>a</sup>, Claudio Gariazzo<sup>e</sup>, Alexei Lyapustin<sup>f</sup>, Paola Michelozzi<sup>a</sup>, Matteo Renzi<sup>a</sup>, Matteo Scortichini<sup>a</sup>, Alexandra Shtein<sup>g</sup>, Giovanni Viegi<sup>h</sup>, Itai Kloog<sup>g</sup>, Joel Schwartz<sup>i</sup>

<sup>a</sup> Department of Epidemiology, Lazio Regional Health Service/ASL Roma 1, Via C. Colombo 112, 00147 Rome, Italy

<sup>b</sup> Karolinska Institutet, Institute of Environmental Medicine, Stockholm, Sweden

<sup>c</sup> Swiss Tropical and Public Health Institute, Basel, Switzerland

<sup>d</sup> University of Basel, Basel, Switzerland

<sup>e</sup> INAIL, Department of Occupational & Environmental Medicine, Monteporzio Catone, Italy

<sup>f</sup> National Aeronautics and Space Administration (NASA) Goddard Space Flight Center (GSFC), Greenbelt, MD, USA

<sup>g</sup> Department of Geography and Environmental Development, Ben-Gurion University of the Negev, Beer Sheva, Israel

<sup>h</sup> Institute of Biomedicine and Molecular Immunology “Alberto Monroy”, National Research Council, Palermo, Italy

<sup>i</sup> Department of Environmental Health, Harvard T. H. Chan School of Public Health, Cambridge, MA, USA

## ARTICLE INFO

Handling Editor: Xavier Querol

### Keywords:

Aerosol optical depth  
Exposure assessment  
Machine learning  
Particulate matter  
Random forest  
Satellite

## ABSTRACT

Particulate matter (PM) air pollution is one of the major causes of death worldwide, with demonstrated adverse effects from both short-term and long-term exposure. Most of the epidemiological studies have been conducted in cities because of the lack of reliable spatiotemporal estimates of particles exposure in nonurban settings. The objective of this study is to estimate daily PM<sub>10</sub> (PM < 10 µm), fine (PM < 2.5 µm, PM<sub>2.5</sub>) and coarse particles (PM between 2.5 and 10 µm, PM<sub>2.5–10</sub>) at 1-km<sup>2</sup> grid for 2013–2015 using a machine learning approach, the Random Forest (RF). Separate RF models were defined to: predict PM<sub>2.5</sub> and PM<sub>2.5–10</sub> concentrations in monitors where only PM<sub>10</sub> data were available (stage 1); impute missing satellite Aerosol Optical Depth (AOD) data using estimates from atmospheric ensemble models (stage 2); establish a relationship between measured PM and satellite, land use and meteorological parameters (stage 3); predict stage 3 model over each 1-km<sup>2</sup> grid cell of Italy (stage 4); and improve stage 3 predictions by using small-scale predictors computed at the monitor locations or within a small buffer (stage 5). Our models were able to capture most of PM variability, with mean cross-validation (CV) R<sup>2</sup> of 0.75 and 0.80 (stage 3) and 0.84 and 0.86 (stage 5) for PM<sub>10</sub> and PM<sub>2.5</sub>, respectively. Model fitting was less optimal for PM<sub>2.5–10</sub>, in summer months and in southern Italy. Finally, predictions were equally good in capturing annual and daily PM variability, therefore they can be used as reliable exposure estimates for investigating long-term and short-term health effects.

## 1. Introduction

Air pollution, especially particulate matter (PM), is one of the major causes of death. Recently, the World Health Organization estimated around 4.2 million of deaths attributable to air pollution exposure worldwide (WHO, 2018). Similarly, the latest update of the Global Burden of Diseases study ranked PM as the sixth leading cause of death out of a list of 84 risk factors, being responsible for over 4 million deaths in 2016 (GBD 2016 Risk Factors Collaborators, 2017).

During the last decades, many epidemiological studies reported consistent health effects of PM from both short-term (i.e. daily

variability) and long-term (i.e. annual averages) exposures, however studies have been historically conducted in major cities, where monitoring networks are more dense and allow measurements and models of spatiotemporal PM variability with more accuracy (Atkinson et al., 2014; Beelen et al., 2014; Brook et al., 2010; Cesaroni et al., 2013; Raaschou-Nielsen et al., 2013; Rückerl et al., 2011; Samoli et al., 2013; Stafoggia et al., 2013). However, it is important to better characterize air pollution distribution and its health effects also in smaller cities, sub-urban and rural areas, where a large fraction of the population lives and which might display higher baseline risks because of less access to healthcare facilities or more deprived socio-economic conditions (Bravo

\* Corresponding author at: Department of Epidemiology of the Lazio Regional Health Service/ASL Roma 1, Via C. Colombo 112, 00147 Rome, Italy.

E-mail address: [m.stafoggia@deplazio.it](mailto:m.stafoggia@deplazio.it) (M. Stafoggia).

<https://doi.org/10.1016/j.envint.2019.01.016>

Received 16 November 2018; Received in revised form 4 January 2019; Accepted 6 January 2019

Available online 14 January 2019

0160-4120/ © 2019 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

et al., 2017; Matz et al., 2015).

In the last ten years there has been an abundance of studies using satellite measurements of Aerosol Optical Depth (AOD), sometimes referred to as Aerosol Optical Thickness, to help predicting ground-level PM concentrations over places or periods with no measurements (de Hoogh et al., 2018; Di et al., 2016; Kloog et al., 2012; Lee et al., 2015; Stafoggia et al., 2016). AOD quantifies the amount of light absorbed or scattered by suspended particles, therefore it represents a relevant parameter to predict PM variability, though an imperfect one, being a columnar estimate while PM concentrations are measured at the ground level. One of the main limitations of AOD is that it is often missing due to cloud coverage, snow or water glint contamination, satellite calibration maneuvers or lost data transmission, which has induced investigators to fill in large gaps in PM predictions by use of different approaches such as kriging, spatiotemporal interpolation or geographic weighted regression. Freely available AOD estimates from atmospheric ensemble models might represent a promising source of data in order to pre-impute satellite-based AOD missing values before PM calibration, however applications in this field are scarce (Zhu et al., 2017).

We have previously developed a spatiotemporal model aimed at predicting daily PM<sub>10</sub> for each 1-km<sup>2</sup> of Italy for the years 2006–2012 using satellite AOD, land use variables and meteorology. Specifically, we have calibrated observed PM<sub>10</sub> concentrations to AOD data using a multivariate linear mixed model, with random intercepts by day and slopes by AOD, and fixed effects for the other spatial and spatiotemporal parameters. The model was able to predict 65% of PM variability on held-out monitors, and displayed negligible bias (Stafoggia et al., 2016).

Machine learning methods have recently become an integral part of modern research, as they offer flexible and automated procedures for the prediction of a target variable based on past observations, unraveling at the same time underlying patterns in data and dealing with complex interactions among predictive variables (Liaw and Wiener, 2002). Among the many different machine learning approaches available, random forests have several advantages compared to other machine learning methods, including: the existence of user-friendly open-source R libraries, among which *ranger* is designed to efficiently handle big data; the simplicity of the method, which requires the selection of only three parameters, *mtry* (the number of variables in the random subset at each node), *num.trees* (the number of trees in the forest) and *min.node.size* (which governs the depth of each regression tree); the robustness of the model to parameter specifications; and the ability of the model (inherent to all decision tree-based designs) to handle non-linearity and high-order interactions among predictive variables.

Random forest methods have been applied to estimate pollutant exposure at large spatial scale. Chen et al. (2018a, 2018b) used the random forest approach to estimate PM<sub>10</sub> and PM<sub>2.5</sub> over China during 2005–2016 providing AOD, meteorology and land use information as predictors. Similarly, Araki et al. (2018) developed a spatiotemporal land use random forest model for estimating metropolitan NO<sub>2</sub> exposure in Japan and Hu et al. (2017) applied random Forests for PM predictions in US.

The ongoing BEEP (Big data in Environmental and occupational Epidemiology) project aims to collect, link and analyze a large amount of data coming from different sources to support exposure assessment and environmental epidemiology studies in Italy. In the frame of the BEEP project, we developed a five-stage modelling strategy based on random forests to impute missing AOD data using atmospheric ensemble models, and to predict PM<sub>10</sub>, PM<sub>2.5</sub> and PM<sub>2.5–10</sub> daily concentrations at 1-km<sup>2</sup> spatial resolution across Italy for the period 2013–2015. As secondary aims, we also provide updated estimates of PM<sub>10</sub> for the period 2006–2012.

## 2. Materials and methods

### 2.1. Study domain

Italy is a boot-shaped peninsula located in southern Europe. It is characterized by diverse geo-climatic areas, with two major mountain ranges (Alps and Apennines), one large plain (the Po valley), a long coastal line and many medium-sized urban areas (46 municipalities above 100,000 inhabitants, 99 between 50,000 and 100,000, 165 between 30,000 and 50,000). Big metropolitan areas are also located along the territory with population over 500,000 inhabitants. The country's total area is 307,635 km<sup>2</sup>. For the aims of the project, we divided the Italy spatial domain into 1-km<sup>2</sup> grid cells, as previously reported (Badaloni et al., 2018).

### 2.2. Data

#### 2.2.1. PM monitored data

We obtained daily data on 24-hour mean PM<sub>10</sub> and PM<sub>2.5</sub> concentrations over the period 2006–2015 from all the available monitoring sites provided by the Italian Institute for Environmental Protection and Research (ISPRA). There were 198, 221 and 229 stations measuring both PM<sub>10</sub> and PM<sub>2.5</sub> during 2013, 2014 and 2015 respectively, while 308, 298 and 295 stations measured only PM<sub>10</sub> concentrations. In each monitor where both PM fractions were available, the coarse fraction (PM<sub>2.5–10</sub>) was derived by subtracting PM<sub>2.5</sub> from PM<sub>10</sub>. Since the availability of PM<sub>2.5</sub> monitors before 2013 was very limited, we restricted our models for PM<sub>2.5</sub> and PM<sub>2.5–10</sub> to 2013–2015.

#### 2.2.2. Aerosol Optical Depth (AOD) data

Aerosol Optical Depth (AOD) is a satellite parameter measuring the degree to which suspended particles affect the transmission of light by absorption or scattering. Therefore it is an indirect measure of the particles present in the column of air on a given time. Recently NASA has developed an algorithm, the Multi-Angle Implementation of Atmospheric Correction (MAIAC), which provides better quality AOD data at 1-km<sup>2</sup> spatial resolution compared with the standard MODIS products (Lyapustin et al., 2018). In this analysis, as in our previous application (Stafoggia et al., 2016), we used MAIAC AOD based on collection 6 MODIS Aqua L1B data for the years 2006–2015.

MAIAC AOD data can be unavailable on a large sample of grid cells and days because of cloud coverage, water/snow glint reflectance and satellite calibration. Therefore, we downloaded modelled AOD data from the Monitoring Atmospheric Composition and Climate - Interim Implementation (MACC-II) project, developed within the Copernicus Atmosphere Monitoring Service (CAMS) and available from the European Centre for Medium-Range Weather Forecasts (ECMWF) website (MACC-II Collaborative Group, 2014). CAMS provides predicted total AOD as the sum of five types of tropospheric aerosols: sea salt, dust, organic and black carbon, and sulfates. Three-hour AOD estimates at five different wavelengths (469 nm, 550 nm, 670 nm, 865 nm and 1240 nm) for all days in 2006–2015 were downloaded at the maximum spatial resolution available in ECMWF, equal to 0.125° × 0.125° (approximately 10 × 10-km<sup>2</sup>).

#### 2.2.3. Meteorological parameters

Meteorological parameters (daily mean air temperature, sea-level barometric pressure, precipitations, relative humidity, wind speed and direction) and planetary boundary layer height were retrieved by the ERA-Interim reanalysis project (Dee et al., 2011), the latest global atmospheric reanalysis produced by the ECMWF. Data have been downloaded at the spatial resolution of 0.125° × 0.125° for the hours 0.00 and 12.00 for each day in 2006–2015.

#### 2.2.4. Other spatiotemporal data

We collected monthly estimates of Normalized Difference

**Table 1**  
Description of the spatial variables.

Variable	Description	Source	Spatial resolution
Domain	307,635 $1 \times 1\text{-km}^2$ grid cells	–	1 km
Administrative areas	Regions, Provinces, Municipalities	ISTAT	Polygons
Geo-climatic zones	Alpine ridge (zone 1), Po valley (zone 2), high Adriatic (zone 3), Appennine (zone 4), high Tyrrenum (zone 5), mid Tyrrenum (zone 6), low Adriatic and Ionium (zone 7), low Tyrrenum and Sicily (zone 8), Sardinia (zone 9)	ISPRA	9 macro-areas
Population	Resident population from census October 2011	ISTAT	Census blocks
Corine land cover	Land cover characteristics	EEA	~100 m
Imperviousness surface areas	An indicator of the spatial distribution of artificial areas. Examples include housing areas, traffic areas (airports, harbors, railway yards, parking lots), roads, industrial and commercial areas, construction sites, etc.	EEA - CLMS	~20 m
Light at night	Satellite-based nighttime imagery	VIIRS - DNB	750 m
Elevation	European Digital Elevation Model EU-DEM	EEA - CLMS	~30 m
Roads	Road density (meters within the cell) and proximity (distance between centroid and closest road) by road type: highway, major, secondary, local	TeleAtlas TomTom network	Lines
Industrial emissions	PM <sub>10</sub> , SO <sub>2</sub> , NO <sub>2</sub> , CO and NH <sub>3</sub> emissions (year 2010) from 743 major industrial plants; proximity from the closest point	ISPRA	Points

Abbreviations: ISTAT (Italian Institute of Statistics), ISPRA (Italian National Institute for Environmental Protection and Research), EEA (European Environment Agency), CLMS (Copernicus Land Monitoring Service), VIIRS-DNB (Visible Infrared Imaging Radiometer Suite-Day/Night Band).

Vegetation Index (NDVI) from the publicly available MODIS NDVI product (MOD13A3) at  $1\text{-km}^2$  spatial resolution. Desert dust advection days were identified across the whole country using a combination of atmospheric tools, and a 0/1 indicator was defined for each grid cell and day based on the absence/presence of a desert advection episode. See [Pey et al. \(2013\)](#) for further details.

### 2.2.5. Spatial data

We computed a number of spatial predictors at the grid cell level, e.g. variables changing from cell to cell but assumed to be fixed over time. These are summarized in [Table 1](#) and include: a) geo-climatic zones, as defined by ISPRA; b) resident population, based on the National Census 2011; c) point emission sources, provided by ISPRA and expressed as tons/year of five pollutants (PM<sub>10</sub>, SO<sub>2</sub>, NO<sub>2</sub>, CO, NH<sub>3</sub>) emitted in 2010 by 743 industrial plants distributed across the country; d) total emissions (emitted from both point and areal sources on 2010) of the five pollutants for each of the 110 Italian provinces; e) mean elevation, obtained from the Copernicus Land Monitoring Service (CLMS) - European Digital Elevation Model (EU-DEM), at 30 m spatial resolution; f) imperviousness surface areas (ISA), derived from CLMS for the year 2012; g) light at night data, collected from the Visible Infrared Imaging Radiometer Suite (VIIRS) Day/Night Band (DNB), year 2015 ([Elvidge et al., 2017](#)); h) land cover data, based on the Corine Land Cover (CLC) database, year 2012 ([EEA, 2013](#)), and defined as percentage of each grid cell covered by nine CLC classes (high/low development, arable/agricultural land, crops, pastures, shrubs, deciduous, evergreen); i) road density data, collected from the TeleAtlas TomTom\_2012 road network, and defined as number of meters within the cell, and distance of the cell centroid from the closest road, for three types of roads: highway, major + secondary, or local, based on CLC Functional Road Classification; j) proximity of each cell centroid to other features: airports, ports, sea, lakes.

### 2.3. Statistical methods

We developed a five-stage machine-learning approach, based on a random forest methodology, aimed at: 1) predicting PM<sub>2.5</sub> and PM<sub>2.5–10</sub> concentrations in monitors where only PM<sub>10</sub> data were available (stage 1), 2) imputing missing MAIAC-AOD data using co-located multi-band CAMS-AOD data (stage 2), 3) calibrating the spatiotemporal PM concentrations with AOD data, meteorological parameters and land-use terms (stage 3), 4) predicting the output of the stage 3 model over all  $1\text{-km}^2$  grid cells of Italy and all days in 2013–2015 (2006–2015 for PM<sub>10</sub>), and 5) improving the stage 3 predictions by using additional information at a finer spatial resolution (monitor coordinates or 150-m buffer),

with the aim of capturing local sources of PM not accounted for by the wider  $1\text{-km}^2$  resolution. Each of the five stages is briefly described below while more details are reported in the online material, appendices A to E. [Fig. 1](#) displays a schematic representation of the five-stage process.

#### 2.3.1. The random forest model

Random forests, in general terms, represent a family of methods that consist in building an *ensemble* (or forest) of decision trees. Different versions of random forests have been proposed in the literature, depending on how data are sampled and decision trees are grown at each iteration ([Breiman, 1994, 2001](#); [Cutler and Zhao, 2001](#); [Geurts et al., 2006](#); [Ho, 1998](#); [Kwok and Carter, 1990](#); [Rodriguez et al., 2006](#)).

In the proper *Random Forest* design ([Breiman, 2001](#), hereafter referred to as RF), each tree is built using a bootstrap sample of the data, and each node of the tree is split according to the best of a subset of randomly chosen predictors ([Liaw and Wiener, 2002](#)). Finally, outputs from each tree are averaged to obtain an *ensemble* prediction of the target variable. The model also provides an estimate of the “importance” of each predictor by quantifying how much prediction error increases when data for that variable is permuted while all others are left unchanged ([Liaw and Wiener, 2002](#)).

In this study we have applied the RF design to each step, as summarized below and described in detail in the online appendices A–E.

#### 2.3.2. Stage 1: predicting PM<sub>2.5</sub> and PM<sub>2.5–10</sub> from PM<sub>10</sub>

The objective of the stage 1 is to estimate the PM<sub>2.5</sub> and PM<sub>2.5–10</sub> data at the monitors by using, as the main predictive variable, daily PM<sub>10</sub> concentrations from co-located stations. The number of monitors measuring PM<sub>10</sub> was 521, 539 and 546 for 2013, 2014 and 2015, respectively. The corresponding figures for PM<sub>2.5</sub> were 198, 221 and 229, distributed across all the 20 Italian regions (before 2013, only 15 regions had PM<sub>2.5</sub> data, preventing the application of this method for the years 2006–2012) ([Fig. A.1](#)). Each monitor is classified according to its location into “traffic”, “industrial” and “background”.

The RF model for this stage is reported in detail in Appendix A. Briefly, for each year in 2013–2015 we defined a RF model where daily PM<sub>2.5</sub> concentrations were the target variable, and co-located PM<sub>10</sub> concentrations were the main predictor. Also, we included monitor location (traffic, industrial or background), month, day of the week and geographical coordinates of the monitor as additional parameters, in order to capture temporal patterns (seasonal and weekly), smooth geographical gradients in PM concentrations distributions and specificity of the PM<sub>2.5</sub>/PM<sub>10</sub> relationship by type of monitoring location in the final predictions.

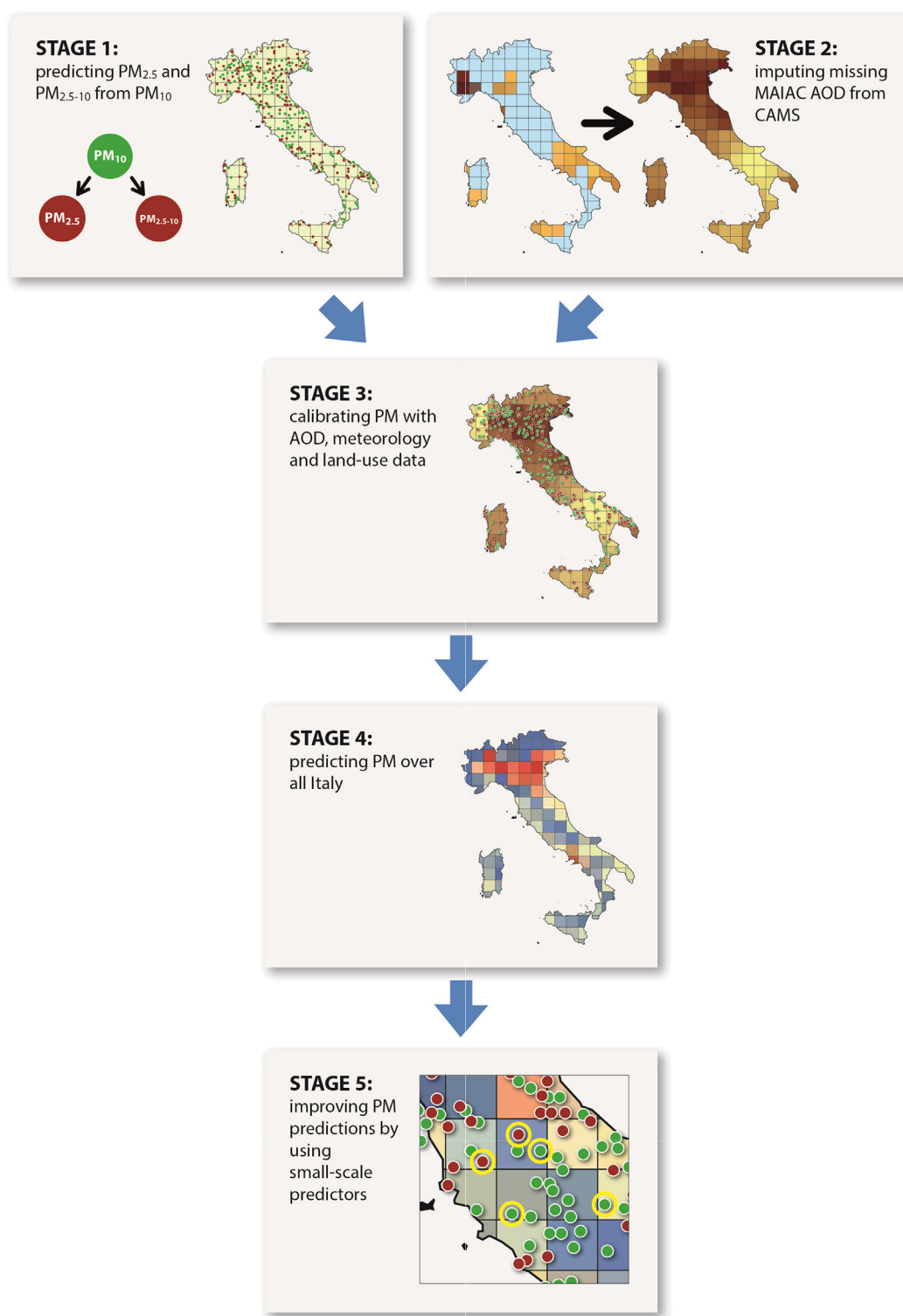


Fig. 1. Graphical representation of the five stage process.

Model fitting was evaluated in two different ways: first, by checking the correlation between observed  $PM_{2.5}$  concentrations and predictions in *out-of-bag* (OOB) samples (i.e. data not included in the bootstrap samples at each iteration of the forest, therefore not contributing to the definition of each decision tree); second, by applying a 10-fold cross-validation approach on the monitoring stations, i.e. randomly splitting the total set of monitors into ten groups, then applying, in turn, the model on nine groups (“training” set) and predicting it to the tenth group (“testing” set), finally checking the correlation between observed

$PM_{2.5}$  concentrations and predictions in held-out monitors. For each of the two approaches we estimated the  $R^2$  (percent of variability of measured  $PM_{2.5}$  captured by predictions), the root mean squared prediction error (RMSPE), and the intercept and slope of the simple linear regression between measured and predicted  $PM_{2.5}$ . The same approach has been applied to  $PM_{2.5-10}$  as the target variable.

#### 2.3.3. Stage 2: imputing missing MAIAC-AOD from CAMS-AOD

As previously mentioned, MAIAC AOD data are often missing. In



**Table 2**

PM<sub>10</sub> and PM<sub>2.5</sub> (observed and predicted in the stage 1) concentrations (µg/m<sup>3</sup>), 2013–2015.

	No. monitors	Mean	SD	Percentiles				
				5	25	50	75	95
PM <sub>10</sub>								
2013	506	25.5	18.1	7.0	14.0	21.0	31.0	61.0
2014	519	24.1	16.9	7.0	13.5	19.9	29.3	56.0
2015	524	26.7	18.2	8.0	15.0	22.0	32.8	62.9
PM <sub>2.5</sub>								
Observed								
2013	198	17.4	14.7	4.3	8.2	13.0	20.8	47.0
2014	221	15.7	12.0	4.4	8.0	12.0	19.0	40.0
2015	229	18.3	14.7	5.0	9.0	14.0	22.1	48.0
Predicted								
2013	506	17.0	14.3	4.7	8.3	12.9	19.9	45.9
2014	519	15.2	11.5	4.6	8.0	12.0	18.0	38.9
2015	524	17.9	14.4	5.0	9.0	13.7	21.3	47.5

Italy, the percentage of missing MAIAC AOD data ranged between 67% in 2011 and 83% in 2014, with larger values in winter and autumn, near the coast and at higher elevations. In our previous study we accounted for such non-random missing patterns by applying inverse-probability weights in the mixed models (Stafoggia et al., 2016). In the present analysis we imputed missing MAIAC AOD values through a RF model using multi-band co-located CAMS AOD values as input variables. Details of the RF model used in this stage are reported in Appendix B. Briefly, for each year in 2006–2015, and separately for the two wavelengths (470 nm and 550 nm) for which MAIAC AOD estimates are provided, we defined a RF model where daily 1-km<sup>2</sup> MAIAC AOD was the target variable and co-located multi-band three-hour estimates of AOD from CAMS were the most relevant input variables. Also, we included day of the year and geographical coordinates of the grid cells centroids as additional parameters, in order to capture residual smooth temporal and spatial patterns in the relationship between MAIAC and CAMS AOD.

Model fitting was evaluated by comparing MAIAC AOD observations and model predictions in the OOB samples. As described in the stage 1 model, we estimated R<sup>2</sup>, RMSPE, intercept and slope as fitting statistics (see Appendix B for further details).

#### 2.3.4. Stage 3: calibrating PM with AOD, meteorology and land-use data

The aim of the stage 3 model is to establish a relationship between daily PM concentrations and AOD, meteorology, and land use data in order to predict PM over locations and days without monitoring stations (stage 4). To this purpose, we developed a RF model having PM (on the log scale) as the target variable (PM<sub>10</sub> from the measurements, PM<sub>2.5</sub> and PM<sub>2.5–10</sub> from the stage 1 estimates) and AOD (as imputed from stage 2 model) and all other spatial and spatiotemporal parameters as potential input variables. We used logarithmic scale to model PM as it displayed a log-normal distribution, and because we wanted to derive non-negative predictions. The details of the model are reported in the online material, Appendix C.

Briefly, for each year and PM metric we modelled log(PM) versus spatiotemporal parameters (including predicted AOD at 470 nm and 550 nm, month, day of the week, meteorological parameters, PBL, NDVI, Saharan dust) + spatial parameters (including cell centroid coordinates, administrative regions, geo-climatic zones, population density, elevation, ISA, light at night, point and areal PM<sub>10</sub> emissions, CLC variables, street density and distance, proximity to airports and sea). RF parameters (*mtry*, *num.trees* and *min.node.size*) have been selected as those which minimized OOB prediction error after a grid search over multiple candidate values.

As the overall aim of the model is to predict PM in locations without PM monitors, model fitting was evaluated comparing observed PM

concentrations with predictions in left-out monitors using 10-fold CV by monitors, as described in stage 1. R<sup>2</sup>, RMSPE, intercept and slope have been computed on the full spatiotemporal predictions and disaggregating between spatial (annual averages) and temporal (difference between daily and annual averages) components, as previously described (Stafoggia et al., 2016).

Finally, we used the (spatial and temporal) intercepts and slopes of the regression between observed and predicted PM in the CV datasets as an estimate of the bias induced by the estimation procedure, and applied 10-fold CV regression calibration on the spatial and temporal components of PM predictions, separately. See Stafoggia et al. (2016) for further details.

#### 2.3.5. Stage 4: predicting PM from stage 3 model over all Italy

We obtained estimates of daily mean PM concentrations for each 1-km<sup>2</sup> grid cell of Italy by applying the stage 3 model fit.

#### 2.3.6. Stage 5: improving PM predictions from stage 3 by using small-scale predictors

The fifth stage is aimed to improve stage 3 PM predictions by capturing local sources of PM variation within grid cell. This is achieved by estimating spatiotemporal predictors at the monitor location (e.g. elevation) or around the monitor (e.g. population and road density within 150-m buffer), and regressing them on the residuals of the stage 3 model. To this aim, we developed an additional RF model, as detailed in the Appendix E. While the output of this stage cannot be applied everywhere (because such small-scale data are not available over all spatial locations of Italy), such data is available in many cohort studies, where it is desirable to predict air pollution concentrations at individual addresses.

All statistical analyses have been performed with the R statistical software, version 3.4.2 (R Development Core Team; <http://R-project.org>). All maps have been produced with ArcGIS software (ESRI. ArcGIS Desktop: Release 10. Redlands, CA: Environmental Systems Research Institute).

### 3. Results

PM monitor locations are displayed in Fig. A.1, Appendix A of the online material, while descriptive statistics of PM<sub>2.5</sub> and PM<sub>10</sub> data (2013–2015) are reported in Table 2 below. For PM<sub>2.5</sub>, both measured and predicted concentrations from stage 1 are displayed. As a consequence, the numbers of monitors for PM<sub>2.5</sub> differs, as they reflect measurements in the upper part, and predictions (which coincide with PM<sub>10</sub> monitors) in the bottom part.

The results of the stage 1 models are displayed in Appendix A, Table A.1 (fitting statistics comparing PM measured concentrations with their predictions in left-out observations and monitors), and Fig. A.2 (scatterplot of observed and 10-fold CV predicted PM concentrations for the year 2015). In summary, PM<sub>2.5</sub> and PM<sub>2.5–10</sub> predictions on testing monitors were unbiased (intercepts close to zero and slopes close to one), with model fitting better for PM<sub>2.5</sub> than PM<sub>2.5–10</sub>, possibly as a consequence of the higher correlation between PM<sub>2.5</sub> and PM<sub>10</sub> compared to PM<sub>2.5–10</sub> and PM<sub>10</sub> (R<sup>2</sup> = 0.93 and 0.60 in 2015, respectively).

Maps with example data for MAIAC and CAMS AOD are presented in Figs. B.1 and B.2 of the online material, Appendix B. Tables B.1 and B.2 in the same Appendix show descriptive statistics of MAIAC and CAMS data, by year and season. Table 3 below reports Pearson correlation coefficients between MAIAC and CAMS AOD data, by year. For the latter, data at 12.00 am have been used, to match as close as possible MAIAC retrievals (as AQUA overpass is at around 1.30 pm daily). Correlations were higher than 0.5 and stable across the years, with highest values at the same wave length 470 nm, as expected.

Table 4 presents the results of Stage 2 models. For each year, fitting statistics from out-of-bag (OOB) samples are reported. Results show very good out-of-sample prediction properties of the RF models, with

**Table 3**

Pearson correlation coefficients between MAIAC (470 nm) and CAMS AOD data (estimated at h12.00, all bands), by year.

Year	AOD				
	469 nm	550 nm	670 nm	865 nm	1240 nm
2006	0.622	0.611	0.592	0.562	0.523
2007	0.635	0.627	0.611	0.581	0.535
2008	0.592	0.589	0.581	0.564	0.535
2009	0.584	0.571	0.549	0.514	0.459
2010	0.619	0.612	0.598	0.570	0.521
2011	0.587	0.581	0.567	0.542	0.500
2012	0.617	0.611	0.598	0.576	0.543
2013	0.601	0.578	0.532	0.449	0.356
2014	0.539	0.523	0.494	0.452	0.416
2015	0.584	0.580	0.499	0.485	0.400

**Table 4**

Results of the stage 2 model. Fitting statistics comparing MAIAC AOD data (at 470 nm, unitless) and predictions, using OOB samples. Table displays  $R^2$  (percent of explained variability), root mean squared prediction error (RMSPE), intercept and slope, by year.

Year	AOD			
	$R^2$	RMSPE	Inter.	Slope
2006	0.957	0.027	0.000	1.004
2007	0.955	0.028	0.000	1.003
2008	0.954	0.027	0.000	1.004
2009	0.946	0.026	0.000	1.004
2010	0.950	0.026	0.000	1.004
2011	0.946	0.026	0.000	1.004
2012	0.942	0.026	0.000	1.003
2013	0.946	0.026	0.000	1.004
2014	0.944	0.025	0.000	1.005
2015	0.949	0.018	0.000	1.003

extremely high  $R^2$  ( $\sim 0.95$ ), negligible mean errors (RMSPE  $\sim 0.02$ ) and no bias (intercepts = 0 and slopes  $\sim 1$ ). There are no differences by years. Fig. B.3 in the appendix displays the maps of predicted MAIAC AOD for two sample days. Finally, Fig. B.4 in the appendix shows the scatterplot of MAIAC-AOD versus CAMS-AOD (left panel) and the scatterplot of MAIAC-AOD versus stage 2 AOD predictions (right panel), for the year 2015.

Table 5 presents the results of the stage 3 model fit for  $PM_{10}$  and  $PM_{2.5}$  (2013–2015). The corresponding results for  $PM_{10}$  (2006–2012) and  $PM_{2.5-10}$  (2013–2015) are presented in the Appendix C, Table C.1. The Stage 3 calibration models for  $PM_{2.5}$  and  $PM_{10}$  all had good out-of-sample predictive performance, with  $R^2 \sim 0.80$  and 0.75 respectively, small prediction errors, negligible bias, and little differences across years.

The most important predictors were spatiotemporal variables (air

**Table 6**

Relative importance (%) of the predictors in the stage 3 model for  $PM_{10}$  and  $PM_{2.5}$  (2013–2015).

Predictor	$PM_{2.5}$			$PM_{10}$		
	2013	2014	2015	2013	2014	2015
Air temperature	13.6	7.2	13.4	7.4	4.4	7.4
PBL (hh 00.00)	9.5	7.3	9.5	9.5	5.7	9.4
Julian day	9.7	10.7	9.2	1.8	8.9	1.8
Barometric pressure	7.9	12.9	7.8	10.2	11.5	9.9
Elevation	7.3	4.7	7.1	9.3	6.7	9.3
PBL (hh 12.00)	6.3	6.2	6.7	8.2	6.0	8.4
Wind (v component)	4.2	4.1	4.0	4.8	5.1	4.6
AOD (470 nm)	2.5	2.5	3.0	2.8	2.8	3.2
AOD (550 nm)	2.5	2.5	2.9	2.8	2.7	3.1
Month	2.9	3.0	2.7	5.1	2.5	5.0
Latitude	2.6	2.9	2.6	3.7	3.7	3.7
Administrative region	2.3	2.0	2.2	0.9	1.5	0.9
Precipitations	1.9	2.5	2.1	3.3	4.0	3.4
Longitude	2.2	1.9	2.1	2.1	1.9	2.1
Wind (u component)	2.0	2.4	2.0	2.6	2.6	2.6
Distance from sea	1.5	1.5	1.4	1.5	1.3	1.5
Resident population	1.4	1.4	1.4	1.5	1.5	1.6
Distance from emission points	1.4	1.9	1.4	1.3	1.4	1.2
Distance from highways	1.3	1.2	1.3	1.2	1.2	1.2
Geoclimatic zone	1.3	1.5	1.3	0.5	1.2	0.5
Density of local streets	1.3	1.9	1.3	1.7	1.6	1.6
$PM_{10}$ emissions from point sources	1.1	2.5	1.2	1.4	1.7	1.3
% Low development	1.1	1.0	1.1	1.1	0.9	1.1
NDVI	1.1	1.3	1.1	1.7	1.8	1.6
$PM_{10}$ emissions from areal sources	1.0	1.1	1.1	1.2	1.0	1.2
Day of week	1.1	1.1	1.1	1.5	1.4	1.5
Distance from airport	1.0	1.2	1.0	1.1	1.3	1.1
% Arable land	0.9	0.9	0.9	0.8	0.8	0.8
Distance from major roads	0.8	0.9	0.8	0.9	0.9	0.8
Light at night	0.8	1.2	0.8	1.0	1.5	1.0
% Deciduous	0.8	0.8	0.8	0.8	0.7	0.7
% Agricultural	0.7	0.8	0.7	0.8	0.9	0.8
Density of major and minor roads	0.7	0.8	0.7	0.9	1.2	0.9
% Shrub	0.6	0.7	0.6	0.8	1.1	0.9
% Crops	0.5	0.5	0.5	0.6	0.5	0.6
Desert dust advection	0.6	1.5	0.5	1.6	3.9	1.5
% High development	0.5	0.5	0.5	0.5	0.6	0.6
% Evergreen	0.4	0.2	0.4	0.4	0.3	0.3
ISA	0.3	0.3	0.3	0.3	0.4	0.3
% Pasture	0.3	0.2	0.3	0.2	0.3	0.2
Density of highways	0.3	0.3	0.3	0.3	0.3	0.3

temperature, PBL, wind components, AOD and Julian day) as they were able to describe PM variability both in space and in time (Table 6). Among the spatial terms, elevation, spatial coordinates and administrative regions showed the highest importance (Table 6).

Poorer fitting was achieved, instead, for  $PM_{2.5-10}$ . There were differences in model fitting by season and geographical area, with model performing worse in summer and southern Italy, while we didn't find differences based on monitor location (Tables C.2 and C.3 of the online

**Table 5**

Results of the stage 3 model for  $PM_{10}$  and  $PM_{2.5}$  (2013–2015). Fitting statistics comparing observed and 10-fold CV predicted PM concentrations, by PM metric and year:  $R^2$  (percent of explained variability), root mean squared prediction error (RMSPE,  $\mu g/m^3$ ), intercept ( $\mu g/m^3$ ) and slope ( $\mu g/m^3$ ), overall and disaggregated by spatial and temporal components.

	Overall				Spatial				Temporal			
	$R^2$	RMSPE	Inter.	Slope	$R^2$	RMSPE	Inter.	Slope	$R^2$	RMSPE	Inter.	Slope
$PM_{10}$												
2013	0.73	9.49	−0.34	1.03	0.62	4.99	1.51	0.95	0.75	8.14	0.00	1.04
2014	0.75	8.40	0.15	1.00	0.61	4.20	1.18	0.96	0.78	7.35	0.00	1.00
2015	0.75	9.05	−0.17	1.01	0.68	4.52	1.22	0.96	0.73	10.18	−0.06	0.99
$PM_{2.5}$												
2013	0.79	6.59	−0.58	1.02	0.76	3.10	−0.03	0.99	0.79	5.88	0.00	1.03
2014	0.78	5.36	−0.54	1.01	0.72	2.49	0.21	0.97	0.79	4.83	0.00	1.02
2015	0.81	6.39	−0.66	1.02	0.79	2.88	−0.10	1.00	0.80	7.05	−0.01	1.01

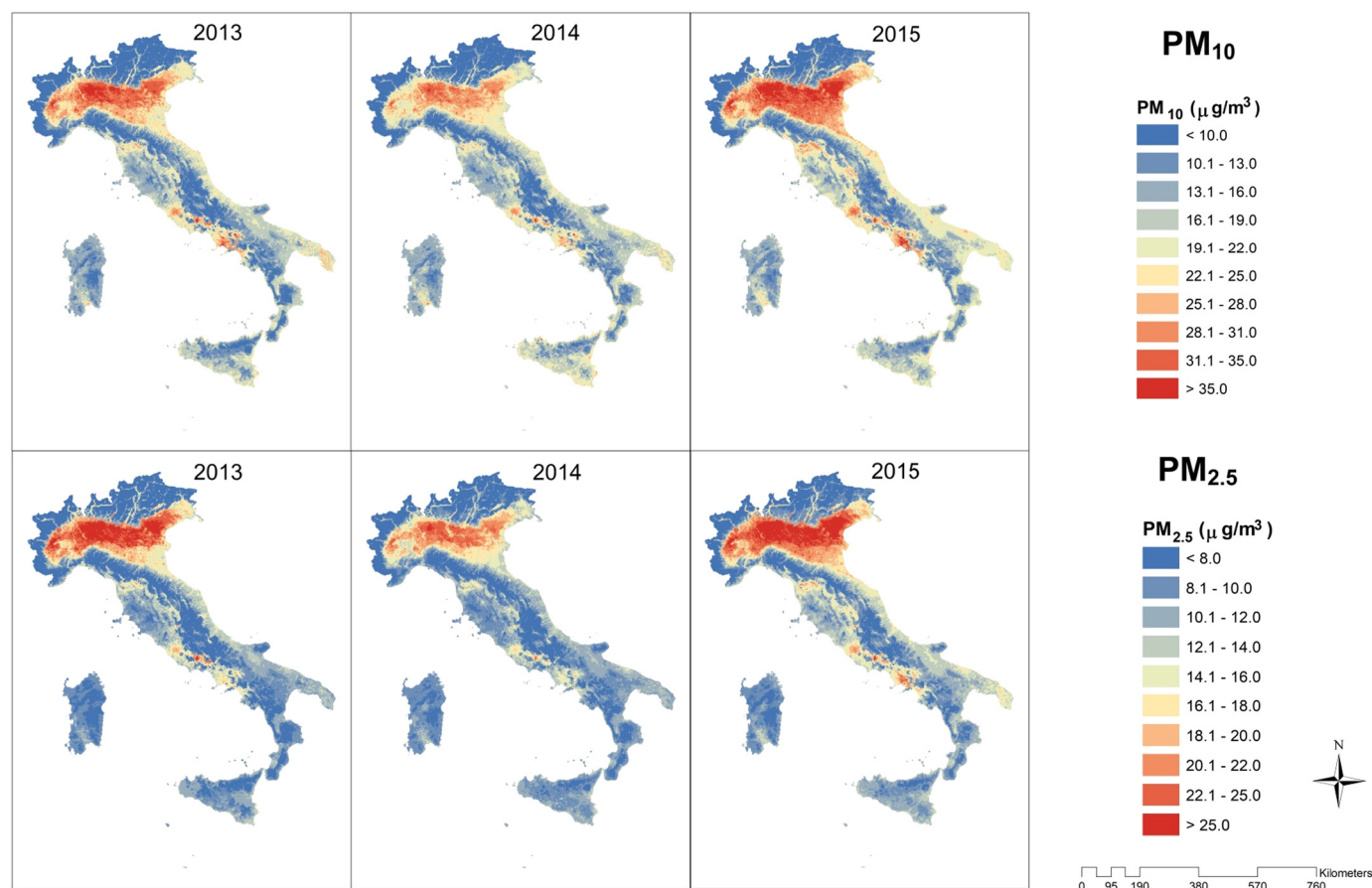


Fig. 2. Predicted  $PM_{10}$  (top) and  $PM_{2.5}$  (bottom) concentrations from stage 4 model: annual means, 2013–2015.

appendix C). The stage 3 performance for  $PM_{10}$  in the years 2006–2012 from the RF approach was slightly improved compared with our previous study (Stafoggia et al., 2016), as displayed in Table C.4 of the Appendix.

The final predictions of  $PM_{10}$  and  $PM_{2.5}$  at the national level for the years 2013–2015 are presented in Fig. 2 and reported as annual average PM concentrations by year. Corresponding results for  $PM_{10}$  (2006–2012) and  $PM_{2.5-10}$  (2013–2015) are presented in the Appendix D, Figs. D.1 and D.2.  $PM_{10}$  and  $PM_{2.5}$  displayed similar spatiotemporal distributions, with concentrations highest in 2015, especially in the major metropolitan areas and in the Po valley.

The stage 5 model (based on small-scale predictors defined around each monitor) substantially improved model fitting. The stage 5 local predictions explained 81%, 87% and 65% of the total variability of  $PM_{10}$ ,  $PM_{2.5}$  and  $PM_{2.5-10}$  respectively, as shown in Table E.1 of the online Appendix E. Predictions were very accurate in capturing both annual average PM concentrations (Fig. 3) and daily means (Fig. 4).

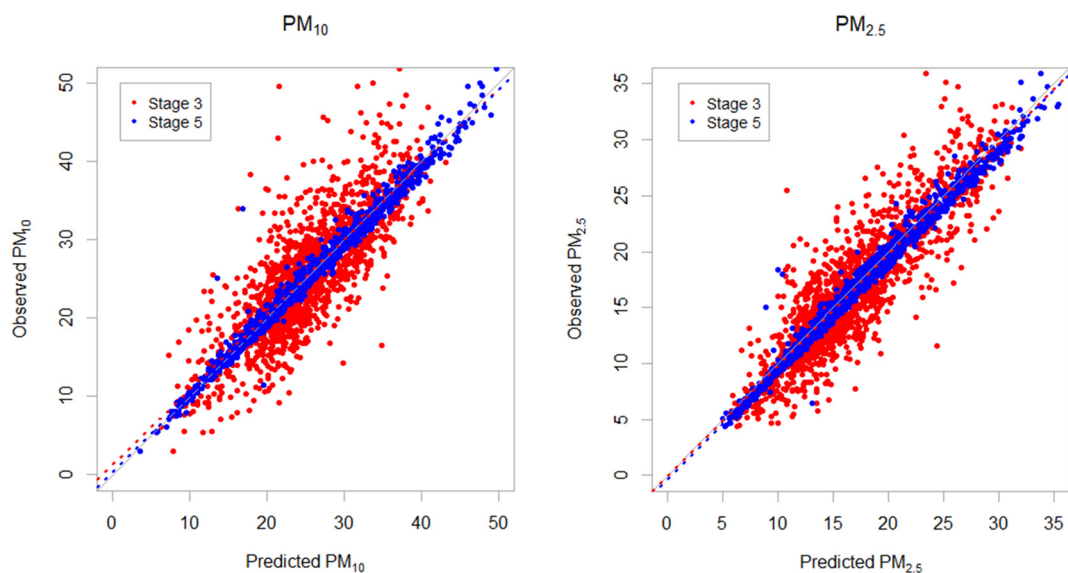
#### 4. Discussion and conclusions

In this study we have developed a five-stage random forest model to impute missing satellite AOD data and predict daily  $PM_{10}$ ,  $PM_{2.5}$  and  $PM_{2.5-10}$  concentrations at fine spatial resolution nationwide. We were able to capture ~75% and 80% of the spatial variability of  $PM_{10}$  and  $PM_{2.5}$  in left-out monitors, with additional 5–10% when small-scale variables were added to predict residuals of the stage 3 model. Model fitting was better in the latest years and in northern Italy, where more monitors are available. Finally, and most importantly, an equally good performance was achieved in predicting day-to-day variability as well as spatial contrasts in annual averages of PM, justifying the use of PM predictions for the analysis of short-term and long-term health effects

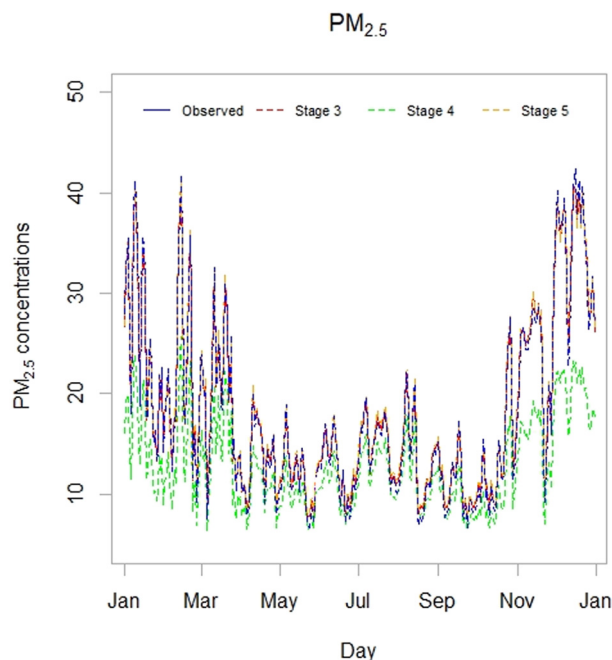
nationwide.

Our models have a number of strength points. First, they improved  $PM_{10}$  predictions compared with our previous mixed-effects model (Stafoggia et al., 2016), both at 1-km<sup>2</sup> level (stage 3) and at the local level (stage 5). In our previous application, we had used linear mixed models with random effects by day. These were a flexible approach to describe temporal patterns, captured by the random intercepts and slopes, but were not as good in describing the complex inter-relationships among the covariates, and the potential non-linearities in the association between them and PM. Adding splines in the models did not help either, because it easily resulted in overfitting the data. Finally, the mixed model could only provide PM predictions in the subset of observations with existing satellite retrievals, while smooth imputation approaches had to be adopted elsewhere (Stafoggia et al., 2016). In the present application we have been able to solve both these problems by applying an additional prediction step for AOD (stage 2) and by using a machine learning method, the RF, explicitly designed for handling complex relationship among predictors without inducing overfitting. In addition, the RF method was robust to parameter specification (number of bootstrap samples, number of predictors used at each split and tree depth) and was computationally efficient, as it allowed to obtain one year of daily predictions (~110 million records) for each PM metric in only a few hours. All this resulted in higher CV-R<sup>2</sup> in both stage 3 and stage 5, with benefits in terms of reduced exposure prediction error for future epidemiological applications.

Second, we were able to fill in missing satellite data by using AOD estimates from atmospheric ensemble models. The results of the stage 2 models were highly stable and accurate, with predictions in OOB samples capturing > 94% of the variability observed in the MAIAC AOD retrievals. This presented the double advantage of allowing us to use all PM data (and not just those intersecting with non-missing



**Fig. 3.**  $PM_{10}$  (left) and  $PM_{2.5}$  (right) average concentrations ( $\mu\text{g}/\text{m}^3$ ) at the 591 monitors available in Italy in 2013–2015: comparison between measured (y axis) and predicted concentrations from stage 3 (red dots) and stage 5 (blue dots) models. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)



**Fig. 4.**  $PM_{2.5}$  concentrations ( $\mu\text{g}/\text{m}^3$ ) in Italy, year 2015: daily averages of  $PM_{2.5}$  measurements (blue line), stage 3 predictions at the  $1\text{-km}^2$  grid cells with monitors (red line), stage 5 predictions at the monitors (orange line), and stage 4 predictions on the whole Italy domain (green line). (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

MAIAC AOD) for stage 3 calibration, and to predict PM across the whole spatiotemporal domain (stage 4), without the need of applying post-hoc smoothing procedures. Third, we were able to predict not only  $PM_{10}$  but also  $PM_{2.5}$  and  $PM_{2.5-10}$ . This is extremely relevant for future epidemiological investigations. Most of the literature shows adverse effects of fine particles (Atkinson et al., 2014; Di et al., 2017; Hoek et al., 2013), as they more easily penetrate and deposit in the lower respiratory tract, possibly translocating into the blood stream and causing adverse effects on the cardiovascular system and in peripheral organs (Brook et al., 2010; Pope 3rd and Dockery, 2006; R  ckel et al., 2011). However,

studies of coarse particles, while fewer than studies of fine particles indicate mortality effects and effects on respiratory health (Behbod et al., 2013; Puett et al., 2009; Stafoggia et al., 2013; Zanobetti and Schwartz, 2009). More studies are clearly needed to clarify their effects. Epidemiological research on  $PM_{2.5}$  health effects in Italy is still scarce because  $PM_{2.5}$  monitors have been installed only recently. The estimates obtained in this study will allow for the first time evaluation of short-term and long-term health effects of fine and coarse particles in Italy, both in the main cities and in smaller cities, sub-urban and rural areas, previously excluded by many epidemiological investigations. Concerning  $PM_{2.5-10}$ , in the last decade the evidence on the potential role of coarse particles as a risk factor to human health has accumulated in the epidemiological literature (Brunekreef and Forsberg, 2005; Keet et al., 2018; WHO, 2013). In particular, coarse particles transported from desert regions frequently impact air quality of (southern) Italy, they have different mineral and chemical composition and might affect health outcomes in a different way (Karanasiou et al., 2012; Perrino et al., 2009). Our spatiotemporal estimates of  $PM_{2.5-10}$  in sub-urban and rural areas with lower  $PM_{2.5}$  concentrations, combined with detection of Saharan dust episodes, will allow investigators to isolate the desert and non-desert contributions to  $PM_{2.5-10}$  and to evaluate their independent health effects.

We also acknowledge some limitations in our approach. Model performance was poorer for  $PM_{2.5-10}$ , in southern Italy and during summer months. Observed data on coarse PM were not from direct measurements but obtained as difference between  $PM_{10}$  and  $PM_{2.5}$ ; therefore they might be affected by two sources of measurement error as well as by the not optimal cutting edge of both observed PM size fractions. As a result, stage 1 model fit for  $PM_{2.5-10}$  was sub-optimal, and this might have worsened the model fit in stage 3. Poorer predictions in southern Italy were expected because of a combination of fewer monitors and less ability of the available predictors to capture the specific PM profile in southern Italy, characterized by large contributions from desert regions, only marginally accounted for in our study. A lower performance of the calibration model in summer months was also found in our previous application (Stafoggia et al., 2016) and deserves further investigation. Another limitation of the study is the hierarchical structure of the models, where outputs from stage 1 ( $PM_{2.5}$  and  $PM_{2.5-10}$  data estimated from co-located  $PM_{10}$ ) and stage 2 (MAIAC AOD imputed from CAMS) serve as inputs for the calibration model in stage 3. This approach prevents a correct quantification of the total



uncertainties of the final PM predictions. On the other hand, the high  $R^2$  and the negligible bias of PM predictions estimated in left-out monitors suggest that, globally, these errors should not be too high and that our model can be exported to locations without monitoring stations with good confidence. It should be acknowledged, however, that the generalization of the calibration model to the entire national domain relies on the assumption that monitor locations are representative of the whole territory, conditional on the geographic covariates. This might not be true, since monitors are oversampled in proximity to traffic sources, residential areas or industrial sites, making estimates more uncertain in remote areas where only few measurements exist. While this might affect the overall layout of our final prediction maps, it is less of a concern from an epidemiological perspective, as these areas are likely underpopulated.

In conclusion, we developed a five stage approach where we merged multiple sources of spatial and temporal data, we predicted satellite AOD from atmospheric ensemble models, and we took full advantage of machine learning methods to obtain finely resolved PM predictions over large spatial and temporal domains. We also applied a local model (stage 5) with the aim of proving the validity of our approach for future epidemiological applications with individual data on residential addresses.

We believe that machine learning methods, in combination with extensive data collection on multiple parameters, can be valid tools for predicting ground level air pollutants concentrations at fine spatial and temporal resolution. While the theory behind machine learning methods is still under development (Jordan and Mitchell, 2015), and more research is required to better characterize all the possible sources of uncertainties inherent to such large estimation processes, we think that our PM predictions will provide novel evidence on the short-term and long-term health effects of fine and coarse particles in Italy.

## Acknowledgments

This paper has been partially funded by the National Institute for Insurance against Accidents at Work, within the project “BEEP” (project code B72F17000180005), and by the Ministry of Health, Italy, within the project “DATISAT” (project code GR-2013-02358899). The funding agencies had no role in the study design; in the collection, analysis and interpretation of data; and in the decision to submit the article for publication. The authors declare no competing financial interests.

## Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.envint.2019.01.016>.

## References

- Araki, S., Shima, M., Yamamoto, K., 2018. Spatiotemporal land use random forest model for estimating metropolitan  $\text{NO}_2$  exposure in Japan. *Sci. Total Environ.* 634, 1269–1277.
- Atkinson, R.W., Kang, S., Anderson, H.R., et al., 2014. Epidemiological time series studies of  $\text{PM}_{2.5}$  and daily mortality and hospital admissions: a systematic review and meta-analysis. *Thorax* 69, 660–665.
- Badaloni, C., Cattani, G., De' Donato, F., et al., 2018. Big data in environmental epidemiology. Satellite and land use data for the estimation of environmental exposures at national level. *Epidemiol. Prev.* 42, 46–59 (Italian).
- Beelen, R., Raaschou-Nielsen, O., Stafoggia, M., et al., 2014. Effects of long-term exposure to air pollution on natural-cause mortality: an analysis of 22 European cohorts within the multicentre ESCAPE project. *Lancet* 383, 785–795.
- Behbod, B., Urch, B., Speck, M., et al., 2013. Endotoxin in concentrated coarse and fine ambient particles induces acute systemic inflammation in controlled human exposures. *Occup. Environ. Med.* 70, 761–767.
- Bravo, M.A., Ebisu, K., Dominici, F., et al., 2017. Airborne fine particles and risk of hospital admissions for understudied populations: effects by urbanicity and short-term cumulative exposures in 708 U.S. counties. *Environ. Health Perspect.* 125, 594–601.
- Breiman, L., 1994. Bagging predictors. Technical report no. 421. Available at: <https://www.stat.berkeley.edu/~breiman/bagging.pdf>, Accessed date: 10 November 2018.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 45, 5–32.
- Brook, R.D., Rajagopalan, S., Pope 3rd, C.A., et al., 2010. American Heart Association Council on Epidemiology and Prevention, Council on the Kidney in Cardiovascular Disease, and Council on Nutrition, Physical Activity and Metabolism. Particulate matter air pollution and cardiovascular disease: an update to the scientific statement from the American Heart Association. *Circulation* 121, 2331–2378.
- Brunekeef, B., Forsberg, B., 2005. Epidemiological evidence of effects of coarse airborne particles on health. *Eur. Respir. J.* 26, 309–318.
- Cesaroni, G., Badaloni, C., Gariazzo, C., et al., 2013. Long-term exposure to urban air pollution and mortality in a cohort of more than a million adults in Rome. *Environ. Health Perspect.* 121, 324–331.
- Chen, G., Li, S., Knibbs, L.D., et al., 2018a. A machine learning method to estimate  $\text{PM}_{2.5}$  concentrations across China with remote sensing, meteorological and land use information. *Sci. Total Environ.* 636, 52–60.
- Chen, G., Wang, Y., Li, S., et al., 2018b. Spatiotemporal patterns of  $\text{PM}_{10}$  concentrations over China during 2005–2016: a satellite-based estimation using the random forests approach. *Environ. Pollut.* 242, 605–613.
- Cutler, A., Zhao, G., 2001. Pert-perfect random tree ensembles. *Comput. Sci. Stat.* 33, 490–497.
- de Hoogh, K., H  ritier, H., Stafoggia, M., et al., 2018. Modelling daily  $\text{PM}_{2.5}$  concentrations at high spatio-temporal resolution across Switzerland. *Environ. Pollut.* 233, 1147–1154.
- Dee, D.P., Uppala, S.M., Simmons, A.J., et al., 2011. The ERA-interim reanalysis: configuration and performance of the data assimilation system. *Q. J. R. Meteorol. Soc.* 137, 553–597.
- Di, Q., Kloog, I., Koutrakis, P., et al., 2016. Assessing  $\text{PM}_{2.5}$  exposures with high spatio-temporal resolution across the continental United States. *Environ. Sci. Technol.* 50, 4712–4721.
- Di, Q., Wang, Y., Zanobetti, A., et al., 2017. Air pollution and mortality in the Medicare population. *N. Engl. J. Med.* 376 (26), 2513–2522.
- EEA (European Environmental Agency), 2013. Corine Land Cover Technical Guide – Addendum 2000. Technical Report No. 40EEA, Copenhagen, Denmark.
- Elvidge, C.D., Baugh, K., Zhizhin, M., et al., 2017. VIIRS night-time lights. *Int. J. Remote Sens.* 38 (21), 5860–5879.
- GBD 2016 Risk Factors Collaborators, 2017. Global, regional, and national comparative risk assessment of 84 behavioural, environmental and occupational, and metabolic risks or clusters of risks, 1990–2016: a systematic analysis for the Global Burden of Disease Study 2016. *Lancet* 390 (10100), 1345–1422.
- Geurts, P., Ernst, D., Wehenkel, L., 2006. Extremely randomized trees. *Mach. Learn.* 63, 3–42.
- Ho, T.K., 1998. The random subspace method for constructing decision forests. *IEEE Trans. Pattern Anal. Mach. Intell.* 20, 832–844.
- Hoek, G., Krishnan, R.M., Beelen, R., et al., 2013. Long-term air pollution exposure and cardio-respiratory mortality: a review. *Environ. Health* 12 (1), 43.
- Hu, X., Belle, J.H., Xia, M., et al., 2017. Estimating  $\text{PM}_{2.5}$  concentrations in the conterminous United States using the random forest approach. *Environ. Sci. Technol.* 51, 6936–6944.
- Jordan, M.I., Mitchell, T.M., 2015. Machine learning: trends, perspectives, and prospects. *Science* 349, 255–260.
- Karanasiou, A., Moreno, N., Moreno, T., et al., 2012. Health effects from Sahara dust episodes in Europe: literature review and research gaps. *Environ. Int.* 47, 107–114.
- Keet, C.A., Keller, J.P., Peng, R.D., 2018. Long-term coarse particulate matter exposure is associated with asthma among children in Medicaid. *Am. J. Respir. Crit. Care Med.* 197, 737–746.
- Kloog, I., Nordio, F., Coull, B.A., et al., 2012. Incorporating local land use regression and satellite aerosol optical depth in a hybrid model of spatiotemporal  $\text{PM}_{2.5}$  exposures in the Mid-Atlantic states. *Environ. Sci. Technol.* 46, 11913–11921.
- Kwok, S.W., Carter, C., 1990. Multiple decision trees. In: *Proceedings of the Fourth Annual Conference on Uncertainty in Artificial Intelligence*. North-Holland Publishing Co, pp. 327–338.
- Lee, M., Kloog, I., Chudnovsky, A., et al., 2015. Spatiotemporal prediction of fine particulate matter using high-resolution satellite images in the Southeastern US 2003–2011. *J. Expo. Sci. Environ. Epidemiol.* 26, 377–384.
- Liaw, A., Wiener, M., 2002. Classification and regression by random forest. *R News* 2, 18–22. <http://CRAN.R-project.org/doc/Rnews/>.
- Lyapustin, A., Wang, Y., Korkin, S., et al., 2018. MODIS collection 6 MAIAC algorithm. *Atmos. Meas. Tech.* 11, 5741–5765.
- MACC-II Collaborative Group, 2014. Final report MACC-II: monitoring atmospheric composition and climate - interim implementation. Available at: [https://atmosphere.copernicus.eu/sites/default/files/repository/MACCII\\_FinalReport\\_0.pdf](https://atmosphere.copernicus.eu/sites/default/files/repository/MACCII_FinalReport_0.pdf), Accessed date: 10 November 2018.
- Matz, C.J., Stieb, D.M., Brion, O., 2015. Urban-rural differences in daily time-activity patterns, occupational activity, and housing characteristics. *Environ. Health* 14, 88.
- Perrino, C., Canepari, S., Catrambone, M., et al., 2009. Influence of natural events on the concentration and composition of atmospheric particulate matter. *Atmos. Environ.* 43, 4766–4779.
- Pey, J., Querol, X., Alastuey, A., et al., 2013. African dust outbreaks over the Mediterranean Basin during 2001–2011:  $\text{PM}_{10}$  concentrations, phenomenology and trends, and its relation with synoptic and mesoscale meteorology. *Atmos. Chem. Phys.* 13, 1395–1410.
- Pope 3rd, C.A., Dockery, D.W., 2006. Health effects of fine particulate air pollution: lines that connect. *J. Air Waste Manage. Assoc.* 56, 709–742 (Review).
- Puett, R.C., Hart, J.E., Yanosky, J.D., et al., 2009. Chronic fine and coarse particulate exposure, mortality, and coronary heart disease in the Nurses' Health Study. *Environ. Health Perspect.* 117, 1697–1701.
- Raaschou-Nielsen, O., Andersen, Z.J., Beelen, R., et al., 2013. Air pollution and lung cancer incidence in 17 European cohorts: prospective analyses from the European

- Study of Cohorts for Air Pollution Effects (ESCAPE). *Lancet Oncol.* 14, 813–822.
- Rodriguez, J.J., Kuncheva, L.I., Alonso, C.J., 2006. Rotation forest: a new classifier ensemble method. *IEEE Trans. Pattern Anal. Mach. Intell.* 28, 1619–1630.
- Rückerl, R., Schneider, A., Breitner, S., et al., 2011. Health effects of particulate air pollution: a review of epidemiological evidence. *Inhal. Toxicol.* 23, 555–592.
- Samoli, E., Stafoggia, M., Rodopoulou, S., et al., 2013. Associations between fine and coarse particles and mortality in Mediterranean cities: results from the MED-PARTICLES project. *Environ. Health Perspect.* 121, 932–938.
- Stafoggia, M., Samoli, E., Alessandrini, E., et al., 2013. Short-term associations between fine and coarse particulate matter and hospitalizations in Southern Europe: results from the MED-PARTICLES project. *Environ. Health Perspect.* 121, 1026–1033.
- Stafoggia, M., Schwartz, J., Badaloni, C., et al., 2016. Estimation of daily PM<sub>10</sub> concentrations in Italy (2006–2012) using finely resolved satellite data, land use variables and meteorology. *Environ. Int.* 99, 234–244.
- WHO (World Health Organization), 2013. Review of Evidence on Health Aspects of Air Pollution – REVIHAAP Project: Final Technical Report. WHO Regional Office for Europe, Copenhagen, Denmark.
- WHO (World Health Organization), 2018. Available at: [https://www.who.int/en/news-room/fact-sheets/detail/ambient-\(outdoor\)-air-quality-and-health](https://www.who.int/en/news-room/fact-sheets/detail/ambient-(outdoor)-air-quality-and-health), Accessed date: 3 January 2019.
- Zanobetti, A., Schwartz, J., 2009. The effect of fine and coarse particulate air pollution on mortality: a national analysis. *Environ. Health Perspect.* 117, 898–903.
- Zhu, J., Xia, X., Wang, J., et al., 2017. Evaluation of aerosol optical depth and aerosol models from VIIRS retrieval algorithms over North China Plain. *Remote Sens.* 9 (5) (pii 432).