

Document downloaded from:

<http://hdl.handle.net/10251/163978>

This paper must be cited as:

Molina-Gomez, NI.; Calderón-Rivera, DS.; Sierra-Parada, R.; Díaz Arévalo, JL.; López Jiménez, PA. (2021). Analysis of incidence of air quality on human health: a case study on the relationship between pollutant concentrations and respiratory diseases in Kennedy, Bogotá. *International Journal of Biometeorology*. 65(1):119-132.  
<https://doi.org/10.1007/s00484-020-01955-4>



The final publication is available at

<https://doi.org/10.1007/s00484-020-01955-4>

Copyright Springer-Verlag

Additional Information



## 33 INTRODUCTION

34

35 In developing countries, air pollution is among the environmental problems of greatest concern. It is a risk factor for  
36 populations' health, which can affect different age groups more severely. Furthermore, growing urbanization has increased  
37 urban density and populations' proximity to pollution sources. Therefore, it is essential to analyze the impact of atmospheric  
38 pollutants on a population's health. In this vein, epidemiological studies and predictive modeling which employ machine  
39 learning (ML) techniques have been carried out.

40

41 Epidemiological studies consist of designing experimental or observational studies. Experimental studies are randomized and  
42 quasi-experimental trials, in which the researcher has a certain degree of control over the variables. Observational studies  
43 include cohort, case-control, cross-sectional and ecological studies (Kestenbaum 2019). In cohort studies, individuals are  
44 classified in sub-groups, according to exposure to a potential cause of sickness, in which the entire evolution of the cohort is  
45 monitored (Lazcano-Ponce et al. 2000). In case-control studies, a comparison is made between the groups in which the event  
46 occurs, and those in which it does not. Cross-sectional studies analyze the frequency of a health event with respect to the  
47 exposure level of the analyzed individuals or group in a given moment (Hernández and Velasco-Mondragón 2000).  
48 Ecological, correlational and exploratory, studies focus on studying groups with an analysis of geographic areas or different  
49 time periods, and are useful in evaluating multiple exposure levels (Borja-Aburto 2000).

50

51 ML techniques were also used to identify the influence of physical and chemical factors in the population's health. ML can  
52 process large volumes of data, as well as linear and non-linear relationships (Ivanov 2018). ML can also perform  
53 classification and regression tasks through decision trees, artificial neural networks (ANN), support vector machines (SVM)  
54 or through ensemble methods, such as random forests (RF) or adaptive boosting (AdaBoost-AdB). From a data set, ML is  
55 able to identify data patterns and predict their behavior (Kuhn and Johnson 2013). ANN were applied to determine the  
56 influence of physical and chemical stressors in hospital admissions for respiratory and cardiac diseases (Kassomenos et al.  
57 2011; Polezer et al. 2018). Generalized boosting models were applied in the same manner for exposure periods before,  
58 during and after forest fires (Reid et al. 2016). Furthermore, Bayesian kernel regression was used to estimate the function of  
59 response doses and to identify the combination of pollutants responsible for adverse health effects (Bobb et al. 2015).  
60 Moreover, statistical tools such as generalized linear regression, multiple linear regression, logistic regression and ML  
61 techniques (RF, SVM, and ANN) were used to forecast atmospheric pollutant levels that may generate a public health risk  
62 (Huang et al. 2018; Ivanov et al. 2018; Kami 2019; Pandey et al. 2013; Weizhen et al. 2014; Zhan et al. 2017).

63

64 The above referenced studies forecasted pollutants' behavior and health effects from information recorded in a database.  
65 However, there is still a lack of forecasting of possible respiratory disease (RD) hotspots based on variables' spatial  
66 distribution and behavior. The spatial distribution of risk factors and their interactions in territories increase interest in  
67 knowing future spatial scenarios of possible health effects. The use of ML and spatial zoning of these factors facilitate  
68 forecasting variables' behavior by identifying hotspots with a territorial approach, which is more specific than national or  
69 capital city focuses. These scenarios are essential for decision-makers so that they are able to implement measures to  
70 mitigate costs related to the treatment of morbidity and mortality.

71

72 With nearly 8.3 million inhabitants and located along the plateau of the eastern range of the Colombian Andes at an  
73 elevation of 2600 meters above sea level (m.a.s.l), Bogotá is one of the most populated cities in Latin America and one of the  
74 cities with the highest recorded levels of atmospheric pollutants, which represent a risk factor for its population. 21.5% of  
75 medical consultations performed for the productive age population (15 – 65 years old) are related to air pollution (García-  
76 Ubaque et al. 2011). Furthermore, changes in NO<sub>2</sub>, SO<sub>2</sub> and PM<sub>2.5</sub> concentration levels in Bogotá were correlated with  
77 statistically significant effects regarding changes in emergency room visits due to RD by children younger than fifteen years  
78 old, while changes in SO<sub>2</sub>, PM<sub>10</sub> and PM<sub>2.5</sub> were related to changes in emergency room visits for circulatory system diseases  
79 in adults older than sixty years (Rodríguez-Villamizar et al. 2018).

80

81 In addition to the above, several studies have been developed in Bogotá aimed at forecasting its air quality. A combined  
82 linear regression model was used to predict air quality (Westerlund et al. 2014). Moreover, ANN were applied to predict  
83 PM<sub>10</sub> and PM<sub>2.5</sub> concentration levels (Franceschi et al. 2018). As a result of the data analyzed, it was determined that  
84 Kennedy is one of the zones in the city with the highest air pollution levels, which also happens to be one of the most  
85 densely populated localities in Bogotá.

86

87 Among the references consulted, there was no study that combined geostatistical tools and ML to identify specific zones in  
88 which cases of RD may occur due to air quality. Therefore, an ecological case study was conducted by applying these tools  
89 in a specific analysis of Kennedy, setting out to determine the influence of meteorological variables and atmospheric  
90 pollutants on the population's health, establishing not only the most relevant variables, but also the zones of greatest interest  
91 in geographic spaces based on the locality's characteristics. Applying geostatistical tools and ML in an environmental health  
92 study on an atmospheric component is innovative. Furthermore, the local scale of the analysis is emphasized, in addition to  
93 data collection in the field being one of the inputs to feed the ML model. This is the first study of this nature developed in  
94 one of the most populated zones of a city such as Bogotá. The approach created in this study can be applied to different  
95 territories, particularly densely populated areas with high air pollution levels. Different municipalities can anticipate  
96 environmental health situations and reduce the cost of RD treatments by applying these tools.

97

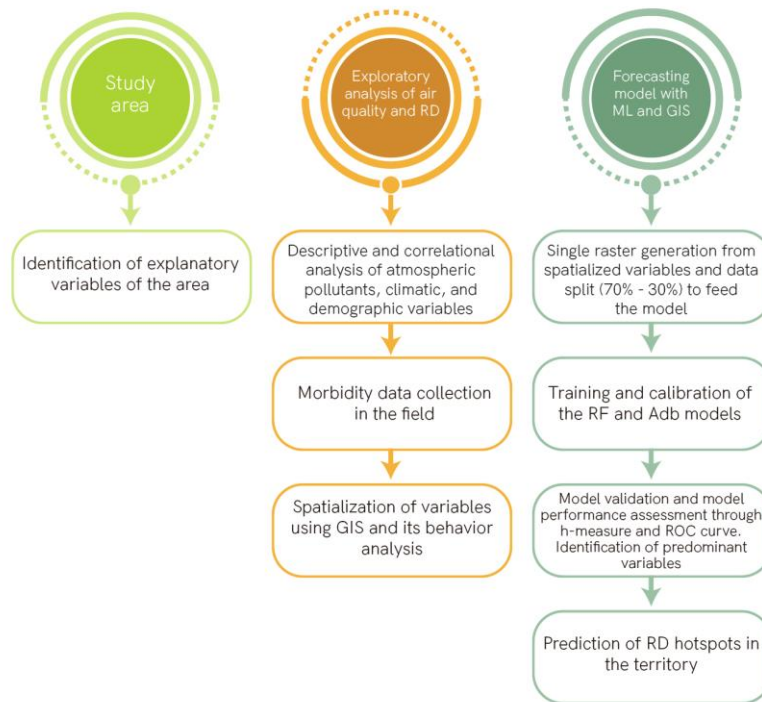
## 98 **MATERIALS AND METHODS**

99

100 This study consists of three sequentially phases which are described below (see Fig. 1), namely: study area analysis,  
101 exploratory analysis of air quality and RD, and the forecasting model with ML techniques and geographical information  
102 system (GIS).

103

104



**Fig. 1 Methodological framework to identify and forecast RD hotspots.**

## STUDY AREA

The study area is Kennedy, located in southwest Bogotá, Colombia. It covers an area of 38.6 km<sup>2</sup>, of which, 1.9% is green space and 10% is protected. Kennedy is located in a transition zone between the eastern plateau and mountains. It is a flat zone which borders five Bogota localities (Puente Aranda, Fontibon, Bosa, Ciudad Bolívar and Tunjuelito), and is where some of the city's main industrial, mining, and commercial activities are located. It also borders Mosquera, Cundinamarca. Wind in Kennedy predominately comes from southwest Bogotá at speeds of 2.2 to 2.5 m/s. The average temperature is 14.6 ± 0.4 °C, with its highest recorded values in 2016 (14.9 ± 0.8 °C). Moderate thermal inversions are common, primarily in the dry months. With respect to precipitation, low values have been recorded in this area of the city, with cumulative averages for the period 2009 – 2016 between 483 and 1018 mm, with a multi-annual average precipitation of nearly 767 mm (SDA 2017).

This locality is made up of twelve zoning planning units (ZPU), which act as territorial units for urban development planning at the zonal level: Américas, Bavaria, Calandaima, Carvajal, Corabastos, Castilla, Gran Britalia, Kennedy Central, Las Margaritas, Patio Bonito, Tintal Norte and Timiza. Four of the above are for residential urban land use,<sup>1</sup> three are for residential use in incomplete urbanization zones,<sup>2</sup> two are in the developing stages,<sup>3</sup> one is the urban center of the locality,<sup>4</sup>

<sup>1</sup> The use changes are occurring in predominately residential sectors with an increase of unplanned territorial occupancy.

<sup>2</sup> Strata 1 & 2 non-consolidated peripheral residential sectors with deficiencies in their infrastructure, accessibility, equipment, and public space.

<sup>3</sup> Under-developed sectors with large unoccupied lots.

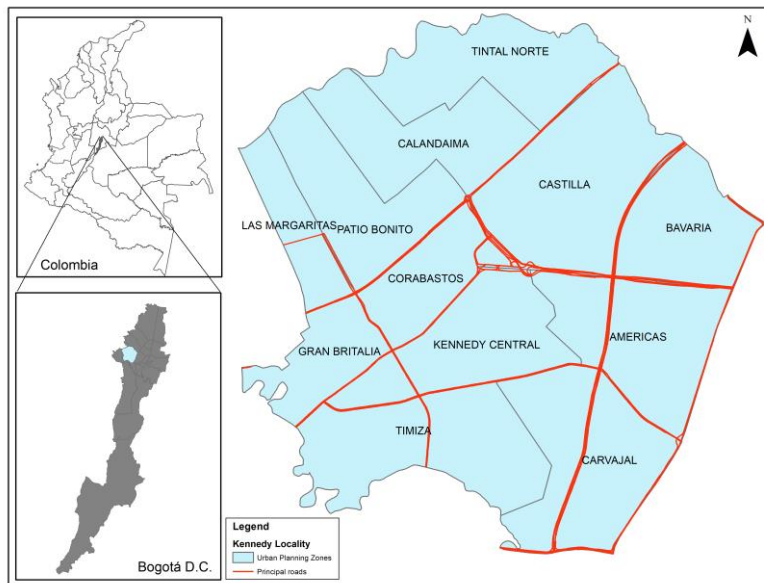
<sup>4</sup> Consolidated sectors that have urban centers with the dominating residential use having been displaced for uses that encourage economic activities.

123 and two are allocated for public use<sup>5</sup> (SDP 2018). There are units registered for industrial development activities in the  
124 Américas, Carvajal and Bavaria ZPUs (Galindo 2013), and approximately 47.2% of households are located near industries  
125 (DANE 2018).

126  
127 Roadways with high vehicular traffic cross the locality, such as Avenida Boyacá, Avenida Ciudad de Cali, and Avenida de  
128 las Américas (see Fig. 2), on which light, cargo and public transportation vehicles represent the main traffic. Predominantly  
129 type 1 and 2 roadways cross and border the locality, as part of the arterial road system. These roads support traffic flows  
130 caused by the inter-urban transport of goods and people. Due to their length and characteristics they support traffic caused by  
131 mass public transportation and connect to the local road network. Type 1 roads are 60 m wide, while type 2 are 40 m wide.  
132 Furthermore, there are local roads within the study area with widths ranging from 4 to 22 m that facilitate entry and local  
133 traffic caused mainly by individual transport vehicles. The locality is also bordered by the Autopista Sur highway and Calle  
134 13, whose road accesses to the city were used by an average of 5300 – 6600 trucks in 2017, with a daily average of 12,000  
135 different types of vehicles going to the CORABASTOS supply center, according to government entities in Bogotá.

136  
137 The predominant buildings in the locality are made of concrete, cinder blocks and bricks, whose heights mainly range from 1  
138 to 5 floors; 67% correspond to buildings up to 3 floors (8.1 m) high, with 18.5% of buildings having 4 and 5 floors (up to  
139 13.5 m), and some buildings are taller than 37 m (DANE 2018).

140  
141 Kennedy has an estimated 1,208,980 inhabitants according to the population census (SDP 2018). In accordance with a 2017  
142 analysis of its population structure, 53.7% of its population are adults, while the early childhood and adolescent population  
143 groups have a smaller representation. Furthermore, the overall population rate in the labor market in Bogotá is approximately  
144 60.8%<sup>6</sup> (SDP 2018).



145  
146 **Fig. 2. Locality of Kennedy and its location in Bogotá, Colombia. The twelve ZPUs that make up the internal**  
147 **distribution of Kennedy along with its main roadways are displayed.**

<sup>5</sup> Large areas allocated to produce urban and metropolitan equipment.

<sup>6</sup> The working age population is 12 years and older in the urban zone, which for Kennedy corresponds to 1,019,894 people.

148

149 **Air quality and its effects on health**

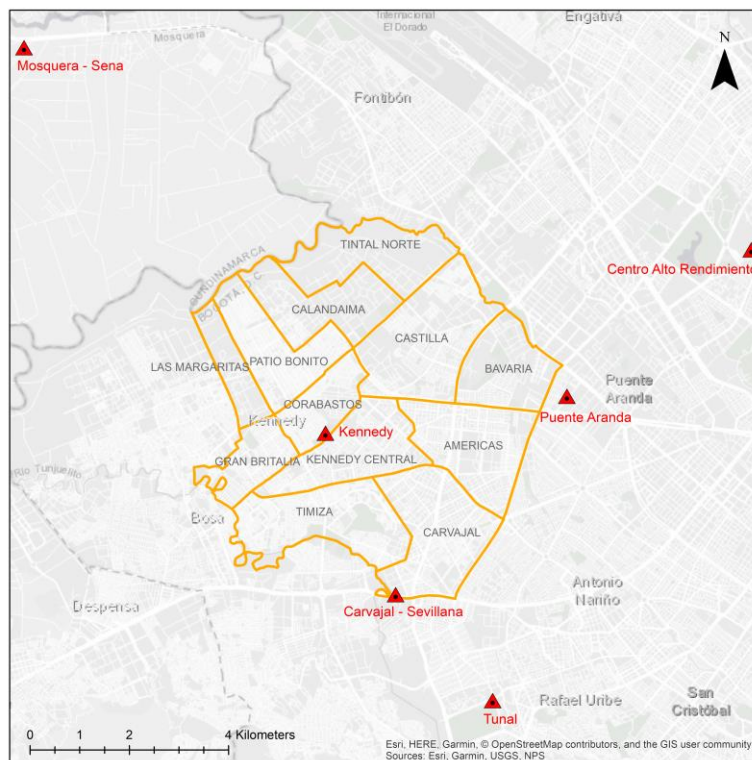
150

151 Bogotá has an Air Quality Monitoring Network (AQMN), which is comprised of thirteen monitoring stations in the city's  
152 urban area, two of which are located in Kennedy (see Fig. 3). The first (Kennedy station) is situated in a residential zone and  
153 monitors PM<sub>10</sub>, PM<sub>2.5</sub>, NO, NO<sub>2</sub>, CO, SO<sub>2</sub>, as well as meteorological variables including humidity, barometric pressure, solar  
154 radiation, temperature, precipitation, wind speed and direction. The second (Carvajal station) is located in a residential zone  
155 with a presence of industrial activity. It is an industrial-traffic station that monitors PM<sub>10</sub>, PM<sub>2.5</sub>, NO, NO<sub>2</sub>, NO<sub>x</sub>, O<sub>3</sub>, CO,  
156 SO<sub>2</sub> and meteorological variables such as precipitation, temperature, wind direction and speed.

157

158 This study analyzed data from the two stations mentioned above, as well as the Tunal, Puente Aranda and Centro de Alto  
159 Rendimiento stations (see Fig. 3), located in Kennedy's influence area, which are the traffic, industrial and background  
160 stations, respectively. These automated stations monitor the same parameters mentioned above as the stations in Kennedy.  
161 Furthermore, the Mosquera-Sena manual station is located in the Bogotá Savanna, and measures PM<sub>10</sub>, SO<sub>2</sub> and NO<sub>2</sub>. The  
162 Centro de Alto Rendimiento station is located in a zone with a low concentration of pollutants, where winds from all  
163 directions converge, and has historically recorded low pollution levels. The Tunal and Puente Aranda stations are located in  
164 zones with high traffic and industrial activities. The Tunal station receives winds from the south, while the winds that hit the  
165 Puente Aranda station come from the west and northwest (SDA 2017).

166



167

168

169

**Fig. 3. Location of the monitoring stations in the study site and the area of influence.**

170 In 2015, the Carvajal and Kennedy stations had the highest concentrations of PM<sub>10</sub> and PM<sub>2.5</sub> in the city. In the first quarter  
171 of 2019, environmental emergencies were declared in areas monitored by these two stations. These emergencies occurred as  
172 the result of a variation in meteorological conditions and the intensification of the temperature inversion phenomena during  
173 the dry season, which generally occurs in the city with moderate effects and breaks atmospheric stability between 7:00 –  
174 9:00 a.m. From 2011 – 2015, neither station met the national regulation standard and recorded yearly concentration averages  
175 that were among the highest for stations that monitor PM<sub>10</sub> and PM<sub>2.5</sub> in the country, which met the temporal coverage  
176 criterion of 75% (IDEAM 2016). NO<sub>2</sub>, O<sub>3</sub>, CO and SO<sub>2</sub> pollutants did not exceed the limits established in the regulation.  
177 However, SO<sub>2</sub> did have higher concentrations in the monitoring stations situated in the locality. The presence of different  
178 types of industrial activities in the city and in neighboring municipalities, as well as road and traffic conditions with different  
179 types of vehicles, all contribute to the increased concentration of atmospheric pollutants in Kennedy. It is important to note  
180 that in Bogotá, prevailing winds from the northeast and southeast displace particulate matter towards the west (Ramírez et al.  
181 2018).

182

### 183 **EXPLORATORY ANALYSIS**

184 Work began on an exploratory analysis of the variables, which in terms of air quality, could impact the population's health in  
185 Kennedy. A descriptive analysis was conducted of the spatial distribution of atmospheric pollutants and meteorological  
186 variables for the period 2009 – 2017, as well as descriptive analyses of individual records from health care providers (RIPS,  
187 as per its Spanish acronym) in Kennedy, reported by the District Health Secretariat (DHS) for the same period, for diseases  
188 associated with air quality, in accordance with Version 10 of the International Classification of Diseases (ICD). Furthermore,  
189 a bivariate Pearson correlation of the continuous variables (pollutants, climatology, and demographic variables) was  
190 conducted for 2016. The data on atmospheric pollutants and meteorological variables was determined based on a weighted  
191 average calculated by the ArcGis 10.5.1 software, using information from the AQMN stations (Carvajal, Kennedy, Puente  
192 Aranda, Centro de Alto Rendimiento and Tunal), as well as the Mosquera-Sena station in Cundinamarca, which are in the  
193 locality and its boarding zones (see Fig. 3).

194

195 The spatial distribution of PM<sub>10</sub>, PM<sub>2.5</sub>, CO, NO<sub>x</sub>, and SO<sub>2</sub> pollutants, as well as the precipitation and temperature  
196 meteorological variables were analyzed via the deterministic method for interpolation called inverse distance weighting  
197 (IDW) interpolation. This is a univariate interpolation method, which is useful in evaluating small study areas. To generate a  
198 predictive surface, the value taken by an unknown point is influenced more by nearby sampled data, than by data from areas  
199 further away (Ly et al. 2011). This method does not consider spatial groupings and has better results when the sampled data  
200 comes from irregularly spaced locations (Li and Heap 2014). This is the case of Bogotá, which has approximately one  
201 station every 23 km<sup>2</sup>. Given that the information for this study comes from irregularly distributed monitoring points (see Fig.  
202 3), this study contemplated examining the influence of the data recorded at the stations concerning its surrounding areas. The  
203 spatial behavioral analysis of the data was performed via the natural grouping data classification method proposed by Jenks  
204 (Jenks 1967).

205

### 206 **FORECASTING MODEL WITH ML AND GEOGRAPHICAL INFORMATION SYSTEM (GIS)**

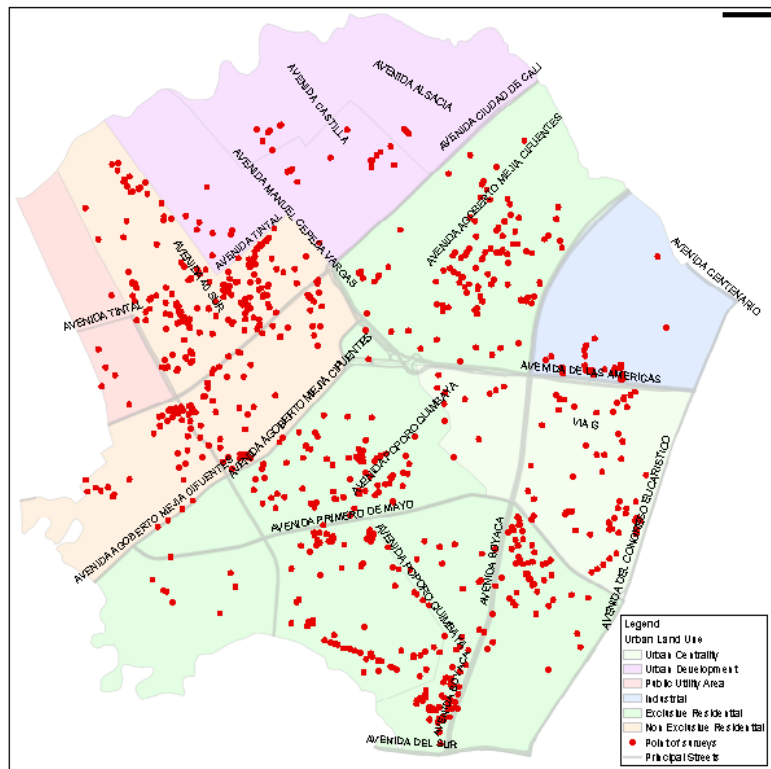
207 Medical consultation records from RIPS do not contain information on the spatial location of health care service users.  
208 Consequently, using data provided by the DHS to spatially identify the zones of the locality with possible RD due to the



209 presence of atmospheric pollutants was not possible. As such, the decision was made to develop the field work by  
210 conducting a survey on health perception, identifying individuals from households in the locality who have been diagnosed  
211 with a RD<sup>7</sup> in 2016. To this end, a georeferenced primary source data collection instrument was applied, which considered  
212 socio-demographic variables and the surveyed person's perception of their health condition. Through a structured  
213 questionnaire, the survey developed which consisted of thirty-one questions, was conducted with households in the locality.  
214 This instrument was applied in twelve ZPUs in 2017 in accordance with the sample size established by the study.

215  
216 The required sample was established to conduct the surveys based on the 2016 number of inhabitants in each ZPU (SDP  
217 2018). The equation for finite populations was used with the following criteria: error: 4%; confidence level: 96%; and  
218 positive and negative variables: 50%. This equation was applied to the general population of the locality yielding a result of  
219 656, which was distributed in accordance with the population proportion of each ZPU. During the development of the study,  
220 the sample size grew to 912 surveys, thus expanding its spatial coverage (see Fig. 4).

221



223

223 **Fig. 4. Field data collection points in the Kennedy ZPUs, primary roadways, and land use typologies**

224

### 225 **Design and development of the forecasting model**

226 Health risks arise from a combination of socio-economic factors, environmental conditions, habitat, and individual behavior.  
227 Geospatial and ML tools were applied to identify the areas of greatest interest related to the population's respiratory health,  
228 in contrast with the presence of pollution sources, pollutant distribution, and the exposed population. The ArcGis 10.5.1

<sup>7</sup> According to the ICD, RD range from rhinopharyngitis, known as the common cold, to respiratory disorders in diseases classified elsewhere (J00-J99).

229 software and the open source software R 3.5.2 (CRAN 2018) were used for this purpose. Information was entered on  
230 atmospheric pollutants (PM<sub>10</sub>, PM<sub>2.5</sub>, CO, NO<sub>x</sub>, SO<sub>2</sub>), meteorological variables, precipitation and temperature previously  
231 determined by the IDW method for 2016, as well as data on population, population density, households' proximity to  
232 roadways (type 1:T1 and type 2:T2), and land use typology. According to Salam et al. (2008) and Li et al. (2011), the  
233 proximity of households, located between 100 m – 1000 m to local and main roads, may increase the risk of RD. This study  
234 considered households' distance to primary and secondary roads (T1 and T2, respectively).

235

236 Geocoded information was also entered and arranged in a GIS of the set of categorical responses from the individuals  
237 surveyed to the question, "In the last year (2016), have you or any of members of your household been diagnosed by a doctor  
238 with a respiratory disease or infection such as asthma, pneumonia or severe lung disease?" A single raster was created with  
239 the resulting information, in which the analyzed variables (continuous and categorical) were overlaid and then used as input  
240 information for the R software. This process provided information on the twelve explanatory variables and the categorical  
241 answer for each point in the locality. RF and the AdB algorithm were the ML tools used, both of which improve the accuracy  
242 of single decision tree classifiers by combining trees grown (Breiman 2001). These tools maintain a bias-variance trade off  
243 through bagging or boosting methods. It is important to note that ANN and SVM are also useful in classification tasks.  
244 However, collinearity of variables is a condition that limits the accuracy and generalization capacity of ANN (Kuhn and  
245 Johnson 2016). Furthermore, the proximity between classes in the geographical space limits the accuracy of SVM. RF and  
246 AdB perform better in those aspects and facilitate the analysis of variables distributed in space, which is useful in the  
247 integrated and spatial analysis of possible health risk factors.

248

249 RF are one of the most accurate bagging methods. RF are a consistent classifier in collecting tree-structured classifiers  
250  $\{h(x, (\Xi)k), k = 1, \dots\}$ , in which  $\{(\Xi)k\}$  are independent random vectors identically distributed, with each tree issuing a single  
251 vote for the most popular class in the  $x$  input (Breiman 2001). For categorical predictions, the voting process selects the class  
252 with the most votes (Kuhn and Johnson 2016). RF can handle large numbers of features (Ivanov et al. 2018) and identify the  
253 most important variables for the model. The precision of RF depends on the strength of the individual classifiers and the  
254 dependence measure between them (Breiman 2001).

255

256 The model used 70% of the data for training and 30% for testing. The partition was performed by randomly considering the  
257 proportionality between affirmative and negative responses. A forecast was created of the areas with the strongest  
258 confluence of affirmative responses to the possibility of RD cases by a majority vote, resulting in the classification that  
259 determined the most influential variables in the model and the distribution of response data according to the conjugate of the  
260 predictor variables in the classification with RF. The model calibration included an iteration of 300 – 1500 trees, with every  
261 100 trees establishing the best combination with the number of variables, according to the accuracy results and the Kappa  
262 index.

263

264 Subsequently, the AdB algorithm, which has no random elements and uses decision trees as the model base, was applied to  
265 create a strong classifier (an ensemble of trees) built from weak classifiers by successively reweighing them (Breiman 2001).  
266 AdB is one of the most widely-used boosting methods in which each classifier focuses on the data that was erroneously  
267 classified by its predecessor, in order to adapt the algorithm and generate better results with each iteration and reduce the

268 generalization error (Schapire and Freund 2012; Rokach and Maimon 2015). In this method, each constructed tree depends  
269 on its predecessor's trees and the prediction come from the most frequent selected class. The samples that are incorrectly  
270 classified in the iteration are given more weight than the samples correctly classified. Therefore, samples that are difficult to  
271 classify are given greater weight until Adb identifies the best model (Kuhn and Johnson 2016). In this study, the same input  
272 parameters for the RF model were used.

273  
274 By using the *Mean Decrease Accuracy* tool, the variables with the greatest influence on the classification error were  
275 determined for each model. Subsequently, forecasts were made with 100% of the spatial behavior data according to possible  
276 RD cases in the locality. The H-measure and the classification error from the receiver operating characteristic (ROC) curve  
277 were used as the performance indicators. The H-measure is a measurement of the loss from erroneous classification  
278 contingent on the relative proportion of the objects belonging to each class (Hand 2009). The ROC curve enables a  
279 comparison of the accuracy and precision of the representing model for each threshold value. This curve is a plot showing all  
280 the sensitivity and specificity pairs resulting from the continuous variation of cutoff over the entire range of observed results  
281 (Altman and Bland 1994). Furthermore, as a function of sensitivity and specificity metrics, the area under the ROC curve  
282 (AUC) is insensitive to disparities in class proportions. A perfect model separates the two classes with sensitivity and  
283 specificity values of 100% (Kuhn and Johnson 2016). Therefore, sensitivity and the specificity metrics of the diagnostic test,  
284 as well as the AUC closest to 1, in the 0.5 – 1.0 interval, represents greater accuracy than the discriminant test (Del Valle  
285 2017). This area establishes the probability that a random person with the disease has a higher measurement value than a  
286 random person without the disease (Altman and Bland 1994).

287  
288 The “Geographic data analysis and modeling” raster and the “Bindings for the 'Geospatial' Data Abstraction Library” rgdal  
289 were the packages used in the R software to read and process the raster images. “Breiman and Cutler's random forests for  
290 classification and regression” were used to design and develop the RF model. The “C\_classification \_A\_nd \_RE\_gression  
291 \_T\_raining” caret was used to determine the most optimal model parameters. The “Visualizing the performance of scoring  
292 classifiers” ROCR was used to calculate the AUC and display the ROC curve. These packages made it possible to adjust the  
293 spatial information to the databases adapted for statistical and predictive processing. A computer with Core i5 8th generation  
294 processor, 8Gb RAM and 1Tb hard disk was used.

## 295 296 **RESULTS**

297  
298 The locality of Kennedy is characterized by its location between primary roadways and the diversity of economic activities  
299 carried out in the same. It has gone through different changes as it is one of the most densely populated localities in Bogotá.  
300 By using the IDW method for the period 2009 – 2017, a decreasing trend was found in the concentration of different  
301 pollutants with values ranging from: (78.72 – 53.11 ug/m<sup>3</sup>) for PM<sub>10</sub>; (35.09 – 24.32 ug/m<sup>3</sup>) for PM<sub>2.5</sub>; (1.2 – 0.73 ug/m<sup>3</sup>) for  
302 CO; (64.15 – 40.46 ppm) for NO<sub>x</sub>; and (8.69 – 2.77 ppb) for SO<sub>2</sub>. The largest values mainly occurred in 2009, which also  
303 had the largest number of consultations associated with RD.

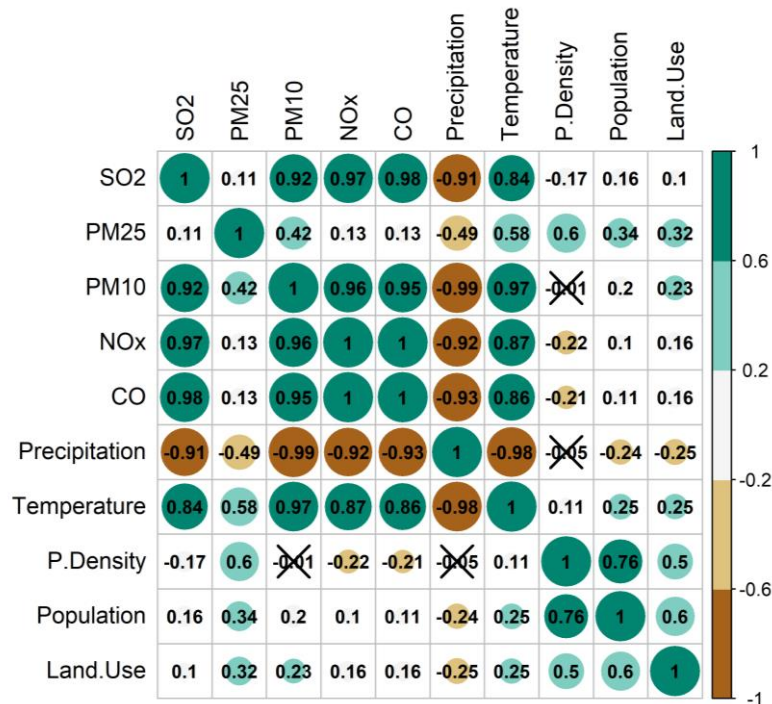
304  
305 The PM<sub>10</sub> and PM<sub>2.5</sub> values surpassed those established by WHO guidelines (PM<sub>10</sub>= 20 µg/m<sup>3</sup>; PM<sub>2.5</sub>= 10 µg/m<sup>3</sup>). According  
306 to the WHO (2006), these are the lowest levels that demonstrate, with more than 95% confidence, that total cardiopulmonary

307 and lung cancer mortality increases in response to prolonged exposure to PM<sub>2.5</sub>. However, values were recorded in 2017 that  
 308 were close to WHO guideline values according to which, the risk of premature mortality is reduced by 6% compared to the  
 309 severe level; (PM<sub>10</sub> =50 µg/m<sup>3</sup>; PM<sub>2.5</sub>= 25 µg/m<sup>3</sup>) (WHO 2006). The SO<sub>2</sub>, CO and NO<sub>x</sub> values indicate a reduction of  
 310 pollutants. SO<sub>2</sub> did not surpass standards (30 ppb) set by the Environmental Protection Agency (EPA). In the case of CO,  
 311 there is no yearly standard, yet the values recorded at the monitoring stations did not, at any time, exceed the Colombian  
 312 standard (5000 ug/m<sup>3</sup>), nor the EPA standard (9000 ug/m<sup>3</sup>) for 8 hours of exposure. NO<sub>x</sub>, which is an unregulated pollutant  
 313 and ozone precursor, decreased by approximately 37% compared to the analyzed periods.

314  
 315 In the period covering 2009 – 2017, after the common cold, chronic obstructive pulmonary disease (COPD), acute  
 316 bronchitis, and unspecified asthma were the most common RDs for which different patients went to consultations.  
 317 Consultations ranged from: (1493 – 5744) for COPD; (2294 – 4736) for acute bronchitis; and (1380 – 2573) for unspecified  
 318 asthma.

319  
 320 **Correlation Analysis**

321 A matrix was created by applying Pearson’s method, which demonstrates high correlation between the 2016 climatic  
 322 variables and atmospheric pollutants analyzed; precipitation and temperature have an inverse relationship (see Fig. 5).  
 323 Moreover, in Fig. 5, the relationships with no significance are marked with an X, that is, their parameter *p-value* is greater  
 324 than 0.05.



325  
 326 **Fig. 5. Correlation matrix of pollutants, meteorological and demographic variables in 2016. The color scale on the**  
 327 **sidebar shows the degree of positive (0 – 1) or negative (0 – -1) correlation between the variables.**  
 328

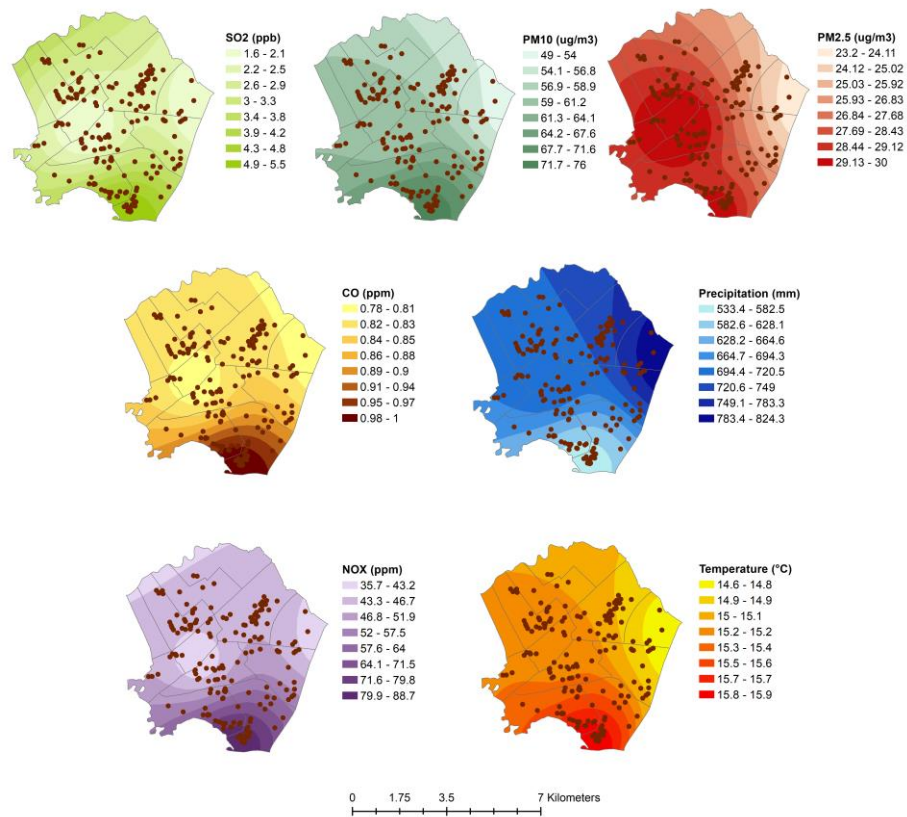
329 **Areas of interest for possible cases of respiratory disease**

330 In 2016, the largest concentration of pollutants analyzed were found in the Carvajal and Timiza ZPUs (see Fig. 6). However,  
 331 the highest concentrations of PM<sub>2.5</sub> were in the areas of the Corabastos, Kennedy Central, Carvajal, Patio Bonito,  
 11

332 Calandaima, Margaritas, Gran Britalia and Timiza ZPUs in the center and western zones of the locality. Temperature had a  
 333 constant behavior (14.98°C), with its lowest values found in the eastern zone of the Bavaria ZPU (14.7°C), which has larger  
 334 precipitation values (769.3 mm) with respect to the rest of the study area (712.42 mm on average). The smallest precipitation  
 335 values were found in the Carvajal and Timiza zones, with 620.6 mm and 637.2 mm, respectively.

336  
 337 In total, 912 surveys were conducted in the twelve ZPUs that make up the locality. The Patio Bonito, Carvajal, Kennedy  
 338 Central and Castilla ZPUs had the largest number of affirmative responses to the questions asked in the field work (see Fig.  
 339 6); 21.4% of the individuals surveyed indicated that a member of their household was diagnosed with a RD, of which 51.8%  
 340 corresponded to the working age population to 60 years old, 23.6% were young people between 5 – 14 years old, and 14.3%  
 341 were people older than 60. Furthermore, it was found that 49.2% of respondents had lived in the study area for more than ten  
 342 years. These indicative figures are comparable with those reported by DHS in 2016, given that in Kennedy nearly 27% of  
 343 RD cases in children under 14 years of age were attended to in emergency rooms, and a prevalence of wheezing was  
 344 reported in 12.6% of adults over 60 years old (SDS 2019).

345



346

347

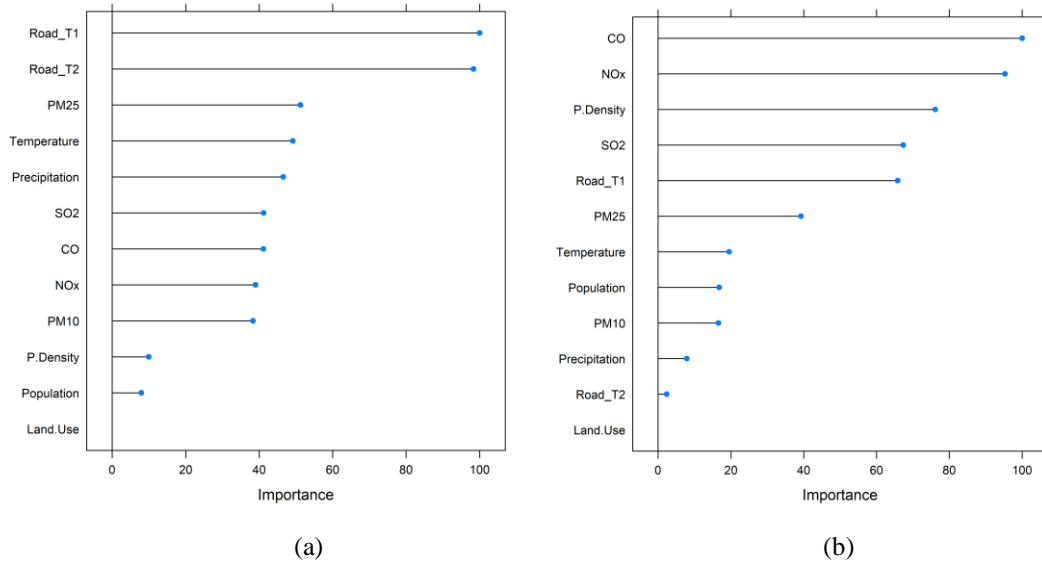
348

**Fig. 6. Behavior of pollutants, meteorological variables, and field work results for 2016**

### 349 Forecasting model with ML and GIS

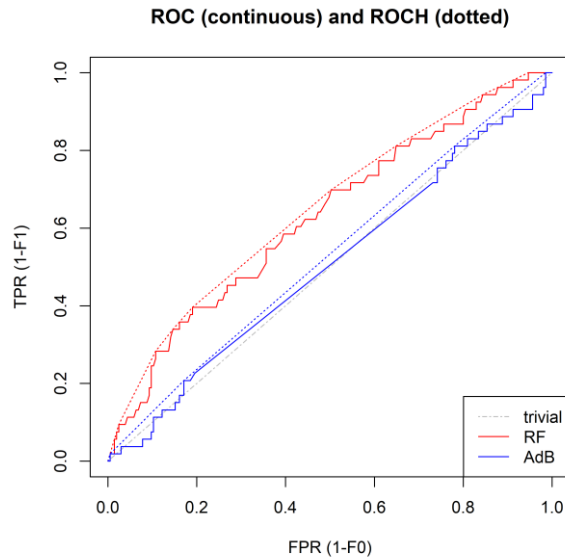
350 Importance matrixes were created (see Fig. 7). The household proximity to roads (T1 and T2 in Fig. 7) variable had the  
 351 strongest influence on the RF model, followed by temperature and PM<sub>2.5</sub>. In the AdB model, household proximity to roads  
 352 was the fifth most important variable. Variable behavior in the model is consistent with respect to the behavior recorded in

353 2016. The population-related variables are the least important in the RF model, while population density (P. Density in Fig.  
 354 7) plays an important role in the AdB model.  
 355



356  
 357  
 358 **Fig. 7. Hierarchization of predictor variables in the (a) RF and (b) AdB models**  
 359

360 With respect to the models' performance, the RF model generated an AUC of 0.63 (see Fig. 8), in which the largest value  
 361 was achieved through a model with 500 trees and 12 variables, which stabilized the error and prevented overfitting, resulting  
 362 in an H measure of 0.10. The AdB model had an AUC of 0.52, for an H measure of 0.018.  
 363

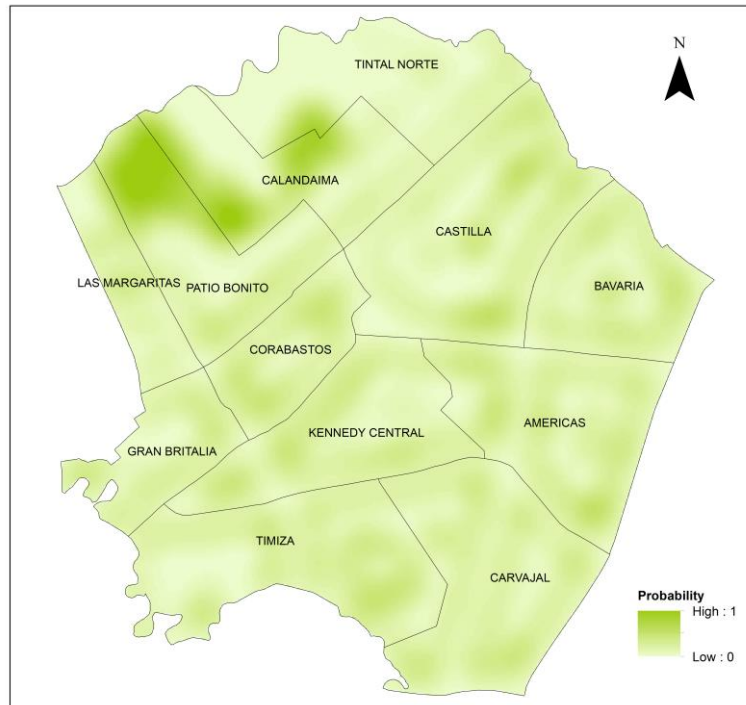


364  
 365 **Fig. 8. ROC curve for RF and AdB**  
 366

367 **Forecasting zones with possible RD events**

368 Based on the behavior of the variables introduced in the RF model, the most relevant zones in the locality related to exposed  
 369 elements and external risk factors are Patio Bonito and Calandaima. Furthermore, considering the confluence of the most

370 important variables in the RF model (proximity to T1 and T2 roads), the behavior of meteorological variables, and pollutants  
371 associated with both road quality and mobile source combustion ( $PM_{2.5}$ , CO,  $NO_x$ ), there are hotspots present in each ZPU,  
372 which, depending on their intensity, enable the occurrence of possible RD cases (see Fig. 9).  
373



374  
375 **Fig. 9. Forecast of zones in which cases of RD could have occurred based on the RF model.**  
376  
377

## 378 DISCUSSION

379  
380 By applying ML tools, it was found that RF outperformed the AdB model for the H-measure, AUC, and accuracy (77.5%  
381 RF; 24.4% AdB). RF had a better odds ratio of 2.25, which reflects a sound diagnosis capacity and greater precision for the  
382 specific classification of possible cases of RD. In the same model, forecasting areas with high densities of possible cases of  
383 RD correlates with proximity to roads,  $PM_{2.5}$ , CO,  $SO_2$  and meteorological variables. The high intensity classification of  
384 areas in Patio Bonito, Calandaima and specific points in Carvajal, Timiza, Kennedy Central and Castilla are supported  
385 primarily by the population's exposure to risk factors, including primary roadways on which different mobile sources of air  
386 pollution transit. However, there are other factors that influence RD events. For example, the low sensitivity of the model  
387 (11%) may be explained by the lack of variables that more accurately describe the behavior of the recorded events. This also  
388 responds to the influence exerted by important variables, as determined by their behavior.

389  
390 It is worth noting that given their spatial behavior, in addition to identifying patterns of pollutant concentrations (Habibi and  
391 Alesheikh 2017), or forecasting individual contaminants (Ivanov et al. 2018), it is necessary to consider the spatial behavior  
392 of combined risk factors and relate this behavior with the exposed population. This fact is sustained in the multiple pollutant  
393 exposure phenomena, which has not been addressed in many studies (Billionnet et al. 2012), as well as the proximity of the  
394 population to pollutants (Mazenq et al. 2017; Yu et al. 2019). This study gathered these experiences and made progress

395 towards an innovative method of identifying zones where possible cases of RD may occur through the combined use of  
396 geostatistical tools and ML.

397

398 In this order of ideas, interpolation via IDW established the behavior of atmospheric pollutants and the meteorological  
399 variables, which were consolidated in the information bases for the correlation matrix and ML models. This aspect is  
400 consistent with prior experiences supported in studies developed by Gorai et al. (2018), Habibi and Alesheikh (2017), and  
401 Sajjadia et al. (2017), in which the transformation of information from observed points to continuous information was  
402 carried out to compare spatial behavior patterns. In different cases, atmospheric pollutant behavior is the primary variable for  
403 the analysis of its effects on a population's health, which is consistent with this study that used ML and found that  $PM_{2.5}$  was  
404 one of the most relevant variables.

405

406 Moreover, when information is collected in the field, conducting surveys has a related bias effect, such as selection bias. To  
407 reduce this effect, the survey was carried out in Kennedy's twelve ZPUs based on the sample, distributed by ZPU and land  
408 use. However, due to security concerns, entering some zones was difficult, which hindered the completion of the total  
409 number of surveys. This was the case in the Las Margaritas and Calandaima ZPUs. Furthermore, in the field data review  
410 process, it was found that due to spatial effects, some regions were not covered. As such, the number of surveys in these  
411 zones was increased to 912 for the final sample size value. Conducting a survey made it possible to identify possible  
412 respiratory system-related morbidity events in a spatial manner. Therefore, related biases may be limited in future uses with  
413 spatial location information that is recorded at health entities, and for security reasons, is not shared with the public.

414

415 Not having complete information of atmospheric pollutants and meteorological variables was one of the study's constraints.  
416 The air quality monitoring and tracking protocol (MAVDT 2010) establishes a minimum data validity standard of 75%. That  
417 is, data whose information is at least 75% complete for the period analyzed is considered valid. Of the data used, an average  
418 of 78% met the temporal validity parameter. The  $NO_x$  and  $SO_2$  data represent 40% of the data that did not meet this  
419 requirement.  $PM_{2.5}$  and CO had values of 38% and 14%, respectively. However, given the need to have information to feed  
420 the ML model and generate a weighted average of pollutant behavior and the meteorological variables for each year, the  
421 recorded information was evaluated in terms of its data trend in an analyzed time series based on the standard deviation of  
422 the variable for the pollutants that did not meet the required validity percentage in a given year and monitoring station. Thus,  
423 the data used exhibited a behavior recorded in its trend, which is largely in line with the required quality standard.

424

425 Detecting hotspots through spatial analysis with geostatistical and ML tools is useful to establish measures to reduce the  
426 vulnerability of people who are exposed to different health risk factors. Moreover, this approach facilitates the identification  
427 of important variables for the model, which is a prioritization tool. Nonetheless, due to different factors that influence  
428 people's health, the model could be strengthened through more available information to refine the characterization of the  
429 study area. This study's approach is useful as a support mechanism for urban planning projects, including the evaluation of  
430 territories' sustainable development performance. This approach could be applied in other fields to identify potential areas of  
431 interest, such as the agriculture sector to identify suitable soils, earth that is ready for the sowing of future crops, or to detect  
432 possible polluted soils due to different activities.

433



434 **CONCLUSIONS**

435

436 - This research developed a tool based on ML that presents the necessary stages to forecast hotspots in which possible RD  
437 cases may occur, based on the behavior of a territory's characterizing variables. Its application in a densely urban area is  
438 useful and replicable as it is a common characteristic in certain territories in developing countries. The micro-territorial  
439 nature of the study is relevant and innovative, as it differs from capital city and country approaches. This approach also  
440 enables researchers to generate useful technical support data for early warnings and contingency plans to mitigate  
441 impacts on air quality and population health, which also influences territories' economies.

442

443 - Using open-source software such as R and spatialization by means of open-source ML codes makes this study an easily  
444 replicable tool. These tools become stronger as more specific and spatialized information becomes available, and their  
445 advantages strengthen environmental health governance by public entities and the academic sector.

446

447 - The level of importance of pollutants such as PM<sub>2.5</sub>, CO, SO<sub>2</sub> and meteorological variables influences the ML model's  
448 behavior. Relevant variables regarding the characteristics of the study area include: high vehicle flow of fossil fuel-  
449 powered automobiles (which explains the level of importance of the PM<sub>2.5</sub>, CO and SO<sub>2</sub> variables); non-standard  
450 operating conditions; deterioration of local roads with the consequent generation of resuspended material; and  
451 residential areas with high population densities that are grouped together, where mixed land uses are integrated with  
452 commercial, industrial and service provision activities. It is necessary to continue carrying out detailed studies on the  
453 exposed population that observe factors such as dose, duration, form of contact, age, sex, diet, personal characteristics,  
454 lifestyle and health condition, in order to determine the relative risk and establish these factors' behavior in the study  
455 area. In this study, 21.4% of the individuals surveyed reported having been diagnosed with respiratory diseases, of  
456 which 14.3% were individuals over 60 years of age, 51.8% are working-aged individuals, and 49.2% of those surveyed  
457 stated that they had lived in the study area for more than 10 years, demonstrating that exposure time is another variable  
458 of interest. These indicative figures include the broad spectrum of respiratory diseases, from the common cold to chronic  
459 and acute respiratory system diseases.

460

461 - Different areas reflect the confluence of risk factors and exposed elements. As such, the RF model established that an  
462 area of great interest could be in the Patio Bonito and Calandaima ZPU's. However, the residential characteristics of the  
463 Timiza, Kennedy Central and Carvajal ZPUs draw attention to the exposed population. RF perform better in terms of a  
464 model driver (AUC: 0.63; H measure: 0.1; accuracy: 77.5%), meaning that the results generated by the RF model are  
465 more accurate than those generated by an AdB model. Similarly, it can be concluded that it is possible to replicate this  
466 model in other areas or municipalities, and its accuracy can be improved by introducing specific data on the location  
467 with the highest exposure of patients attended to in consultations, emergency room visits and hospitalizations related to  
468 RD, as well as information on the explicative variables for the analysis period. The combination of the tools applied in  
469 this study together with a pollutant dispersion model could increase the AUC, as well as the model's classification  
470 metrics.

471

472 - Lastly, it must be stressed that sustainable development refers to an increase in quality of life, through the interaction of  
473 social, environmental, and economic dimensions for equitable, livable, and viable development. A model of these  
474 characteristics becomes a preventive tool, which can contribute to reducing costs by addressing events associated with  
475 air pollution. As a territorial planning component, determining the influence of air pollution on a territory's  
476 sustainability can contribute to implementing policies instituted in the international framework. As air pollution  
477 increases, so does the number of workdays lost, reducing productivity. A better understanding of this phenomenon could  
478 contribute to zonal planning and determining the territorial organization of each zone.

479  
480

#### 481 **Acknowledgments**

482 Many thanks to the members of the Intelligence and Territorial Analysis Group of the Universidad Santo Tomás for their  
483 collaboration in conducting the fieldwork.

484 **References**

- 485 Altman Douglas G, Bland J. Martin (1994) Diagnostic Tests 3: Receiver Operating Characteristic Plots. *BMJ* 309  
486 (6948):188. <https://doi.org/10.1136/bmj.309.6948.188>  
487
- 488 Billionnet C, Sherrill D, Annesi-Maesano I (2012) Estimating the health effects of exposure to multi-pollutant mixture. *Ann*  
489 *Epidemiol* 22:126–141. <https://doi.org/10.1016/J.ANNEPIDEM.2011.11.004>  
490
- 491 Bobb JF, Valeri L, Claus Henn B, et al (2015) Bayesian kernel machine regression for estimating the health effects of multi-  
492 pollutant mixtures. *Biostatistics* 16:493–508. <https://doi.org/10.1093/biostatistics/kxu058>  
493
- 494 Borja-Aburto VH (2000) Ecological studies. *Salud Pública de México* 42:533–538.  
495
- 496 Breiman L (2001) Random Forests. *Machine Learning* 45:5–32. <https://doi.org/10.1023/A:1010933404324>  
497
- 498 CRAN Comprehensive R Archive Network (2018) R-3.5.2 for Windows (32/64 bit). [https://cran.r-](https://cran.r-project.org/bin/windows/base/old/3.5.2/)  
499 [project.org/bin/windows/base/old/3.5.2/](https://cran.r-project.org/bin/windows/base/old/3.5.2/). Accessed 10 June 2019  
500
- 501 DANE National Administrative Department of Statistics (2018) Multi-Purpose Survey -MS 2017. Bogotá, Colombia.  
502
- 503 Del Valle Benavides AR (2017) ROC curves (Receiver-Operating-Characteristic) and their applications. Universidad de  
504 Sevilla.  
505
- 506 DHS (2019) SALUDATA- Health Observatory of Bogota [http://saludata.saludcapital.gov.co/osb/index.php/datos-de-](http://saludata.saludcapital.gov.co/osb/index.php/datos-de-salud/salud-ambiental/consultaagencias14anos/)  
507 [salud/salud-ambiental/consultaagencias14anos/](http://saludata.saludcapital.gov.co/osb/index.php/datos-de-salud/salud-ambiental/consultaagencias14anos/). Accessed 11 April 2019  
508
- 509 Franceschi F, Cobo M, Figueredo M (2018) Discovering relationships and forecasting PM<sub>10</sub> and PM<sub>2.5</sub> concentrations in  
510 Bogotá, Colombia, using Artificial Neural Networks, Principal Component Analysis, and k-means clustering. *Atmos*  
511 *Pollut Res* 9:912-922. <https://doi.org/10.1016/j.apr.2018.02.006>  
512
- 513 Galindo WG (2013) Construction dynamics by use, the locality of Kennedy 2002/2012. Bogotá.  
514
- 515 García-Ubaque JC, García-Ubaque CA, Vaca-Bohórquez ML (2011) Medical consultation in productive age population  
516 related with air pollution levels in Bogota city. *Procedia Environ Sci* 4: 165–169.  
517 <https://doi.org/10.1016/j.proenv.2011.03.020>  
518
- 519 Gorai AK, Tchounwou PB, Biswal S, et al (2018) Spatio-Temporal Variation of Particulate Matter (PM<sub>2.5</sub>) Concentrations  
520 and its health impacts in a mega city, Delhi in India. *Environ Health Insights* 12:1-9.  
521 <https://doi.org/10.1177/1178630218792861>  
522
- 523 Habibi R, Alesheikh AA, Mohammadinia A, et al (2017) An assessment of spatial pattern characterization of air pollution: A  
524 case study of CO and PM<sub>2.5</sub> in Tehran, Iran. *ISPRS Int J Geo-Inf* 6:270. <https://doi.org/10.3390/ijgi6090270>  
525
- 526 Hand DJ (2009) Measuring classifier performance: A coherent alternative to the area under the ROC curve. *Mach Learn*  
527 77:103–123. <https://doi.org/10.1007/s10994-009-5119-5>  
528
- 529 Hernández B, Velasco-Mondragón HE (2000) Cross-sectional surveys. *Salud Pública de México* 42: 447–455  
530
- 531 Huang K, Xiao Q, Meng X, et al (2018) Predicting monthly high-resolution PM<sub>2.5</sub> concentrations with random forest model  
532 in the North China plain. *Environ Pollut* 242:675–683. <https://doi.org/10.1016/j.envpol.2018.07.016>  
533
- 534 IDEAM Institute of Hydrology, Meteorology and Environmental Studies (2016) State of air quality in Colombia, 2011 –  
535 2015 Report. Bogotá D.C.  
536
- 537 Ivanov A, Voynikova D, Stoimenova M, et al (2018) Random forests models of particulate matter PM<sub>10</sub>: A case study, in:  
538 *AIP Conference Proceedings* 2025, 030001. <https://doi.org/10.1063/1.5064879>  
539
- 540 Jenks, George F (1967) The Data Model Concept in Statistical Mapping. *International Yearbook of Cartography* 7: 186-190

541  
542 Kami JA (2019) A random forest partition model for predicting NO<sub>2</sub> concentrations from traffic flow and meteorological  
543 conditions. *Sci Total Environ* 651:475–483. <https://doi.org/10.1016/j.scitotenv.2018.09.196>  
544  
545 Kassomenos P, Petrakis M, Sarigiannis D, et al (2011) Identifying the contribution of physical and chemical stressors to the  
546 daily number of hospital admissions implementing an artificial neural network model. *Air Qual Atmos Health* 4:263–  
547 272. <https://doi.org/10.1007/s11869-011-0139-2>  
548  
549 Kestenbaum B (2019) *Epidemiology and Biostatistics*. Seattle, USA. <https://doi.org/10.1007/978-3-319-96644-1>  
550  
551 Kuhn, Max, Kjell Johnson (2016) *Applied Predictive Modeling*. New York, USA. <https://doi.org/10.1007/978-1-4614-6849-3>  
552 3  
553  
554 Lazcano-Ponce E, Fernández E, Salazar-Martínez E, et al (2000) Cohort studies. Methodology, biases and application. *Salud*  
555 *Pública de México* 42:230–241  
556  
557 Li S, Batterman S, Wasilevich E, et al (2011) Asthma exacerbation and proximity of residence to major roads: a population-  
558 based matched case-control study among the pediatric Medicaid population in Detroit, Michigan. *Environ Health*  
559 10:34. <https://doi.org/10.1186/1476-069X-10-34>  
560  
561 Li Jin, Heap Andrew D (2014) Spatial interpolation methods applied in the environmental sciences: A review. *Environ*  
562 *Modell Software* 53:173-189. <http://dx.doi.org/10.1016/j.envsoft.2013.12.008>  
563  
564 Ly S, Charles C, Degr A (2011) Geostatistical interpolation of daily rainfall at catchment scale: the use of several variogram  
565 models in the Ourthe and Ambleve catchments, Belgium. *Hydrol Earth Syst Sci* 15:2259-2274.  
566 <https://doi.org/10.5194/hess-15-2259-2011>  
567  
568 MAVDT Ministry of Environment, Housing and Territorial Development (2010) Protocol for air quality monitoring. Bogota,  
569 Colombia.  
570  
571 Mazenq J, Dubus J-C, Gaudart J, et al (2017) City housing atmospheric pollutant impact on emergency visit for asthma: A  
572 classification and regression tree approach. *Respir Med* 132:1–8. <https://doi.org/10.1016/j.rmed.2017.09.004>  
573  
574 Pandey G, Zhang B, Jian L (2013) Predicting submicron air pollution indicators: a machine learning approach. *Environ*  
575 *Sci Processes Impacts*. 15:996–1005. <https://doi.org/10.1039/c3em30890a>  
576  
577 Polezer G, Tadano YS, Siqueira HV, et al (2018) Assessing the impact of PM<sub>2.5</sub> on respiratory disease using artificial neural  
578 networks. *Environ Pollut* 235:394-403. <https://doi.org/10.1016/j.envpol.2017.12.111>  
579  
580 Ramírez O, Sánchez de la Campa AM, Amato F, et al (2018) Chemical composition and source apportionment of PM<sub>10</sub> at an  
581 urban background site in a high–altitude Latin American megacity (Bogota, Colombia). *Environ Pollut* 233:142–155.  
582 <https://doi.org/10.1016/j.envpol.2017.10.045>  
583  
584 Reid CE, Jerrett M, Tager IB, et al (2016) Differential respiratory health effects from the 2008 northern California wildfires:  
585 A spatiotemporal approach. *Environ Res* 150:227–235. <https://doi.org/10.1016/J.ENVRES.2016.06.012>  
586  
587 Rodríguez-Villamizar LA, Rojas-Roa NY, Blanco-Becerra LC, et al (2018) Short-Term effects of air pollution on respiratory  
588 and circulatory morbidity in Colombia 2011–2014: A multi-city, time-series analysis. *Int J Environ Res Public Health*  
589 15:2-12. <https://doi.org/10.3390/ijerph15081610>  
590  
591 Rokach, Lior, and Oded Maimon (2015) *Data Mining with Decision Trees: Theory and Applications*. 2nd ed. Singapore:  
592 World Scientific Publishing Co. Pte. Ltd. 5.  
593  
594 Salam MT, Islam T, Gilliland FD (2008) Recent evidence for adverse effects of residential proximity to traffic sources on  
595 asthma. *Curr Opin Pulm Med* 14:3–8. <https://doi.org/10.1097/MCP.0b013e3282f1987a>  
596

597 Sajjadia SA, Zolfagharib G, Adabc H, et al (2017) Measurement and modeling of particulate matter concentrations:  
598 Applying spatial analysis and regression techniques to assess air quality. *MethodsX* 4:372–390.  
599 <https://doi.org/10.1016/j.mex.2017.09.006>  
600

601 Schapire RE, Freund Y (2012) *Boosting: foundations and algorithms*, Adaptive computation and machine learning. MIT  
602 Press, London.  
603

604 SDA District Secretariat for the Environment (2017) *Air quality annual report of Bogota, 2016*. Bogotá, Colombia.  
605

606 SDP District Planning Secretariat (2018) *Monograph 2017 Assessment of the main territorial, infrastructure, demographic*  
607 *and socio-economic aspects of the locality of Kennedy 08*. Bogotá, Colombia.  
608

609 Weizhen H, Zhengqiang L, Yuhuan Z, et al (2014) Using support vector regression to predict PM<sub>10</sub> and PM<sub>2.5</sub>, in: *IOP*  
610 *Conference Series: Earth and Environmental Science*. IOP. <https://doi.org/10.1088/1755-1315/17/1/012268>  
611

612 Westerlund J, Urbain JP, Bonilla J (2014) Application of air quality combination forecasting to Bogota. *Atmos Environ*  
613 89:22-28. <https://doi.org/10.1016/j.atmosenv.2014.02.015>  
614

615 WHO (2006) *WHO Air quality guidelines for particulate matter, ozone, nitrogen dioxide and sulfur dioxide. Global update*  
616 *2005*. Geneva, Switzerland.  
617

618 Yu Y, Yao S, Dong H, et al (2019) Association between short-term exposure to particulate matter air pollution and cause-  
619 specific mortality in Changzhou, China. *Environ Res* 170:7–15. <https://doi.org/10.1016/j.envres.2018.11.041>  
620

621 Zhan Y, Luo Y, Deng X, et al (2017) Spatiotemporal prediction of continuous daily PM<sub>2.5</sub> concentrations across China using  
622 a spatially explicit machine learning algorithm. *Environ Pollut* 233:464-473.  
623 <https://doi.org/10.1016/j.atmosenv.2017.02.023>