

Machine Learning Models Accurately Model Ozone Exposure during Wildfire Events

Gregory L. Watson^{a,*}, Donatello Telesca^a, Colleen E. Reid^b, Gabriele G. Pfister^c, Michael Jerrett^d

^a*Department of Biostatistics, University of California, Los Angeles, California 90024, USA*

^b*Department of Geography, University of Colorado Boulder, Boulder, Colorado 80309, USA*

^c*Atmospheric Chemistry Observations and Modeling Laboratory, National Center for Atmospheric Research, Boulder, Colorado 80301, USA*

^d*Department of Environmental Health Sciences, University of California, Los Angeles, California 90024, USA*

Abstract

Epidemiologists use prediction models to downscale (i.e., interpolate) air pollution exposure where monitoring data is insufficient. This study compares machine learning prediction models for ground-level ozone during wildfires, evaluating the predictive accuracy of ten algorithms on the daily 8-hour maximum average ozone during a 2008 wildfire event in northern California. Models were evaluated using a leave-one-location-out cross-validation (LOLO CV) procedure to account for the spatial and temporal dependence of the data and produce more realistic estimates of prediction error. LOLO CV avoids both the well-known overly optimistic bias of k -fold cross-validation on dependent data and the conservative bias of evaluating prediction error over a coarser spatial resolution via leave- k -locations-out CV. Gradient boosting was the most accurate of the ten machine learning algorithms with the lowest LOLO CV estimated root mean square error (0.228) and the highest LOLO CV \hat{R}^2 (0.677). Random forest was the second best performing algorithm with an LOLO CV \hat{R}^2 of 0.661. The LOLO CV estimates of predictive accuracy were less optimistic than 10-fold CV estimates for all ten models. The difference in estimated accuracy between the 10-fold CV and LOLO CV was greater for more flexible models like

*Corresponding author

Email address: gwatson@ucla.edu (Gregory L. Watson)

gradient boosting and random forest. The order of estimated model accuracy depended on the choice of evaluation metric, indicating that 10-fold CV and LOLO CV may select different models or sets of covariates as optimal, which calls into question the reliability of 10-fold CV for model (or variable) selection. These prediction models are designed for interpolating ozone exposure, and are not suited to inferring the effect of wildfires on ozone or extrapolating to predict ozone in other spatial or temporal domains. This is demonstrated by the inability of the best performing models to accurately predict ozone during 2007 southern California wildfires.

Capsule: Flexible machine learning methods model ozone well during a wildfire. LOLO CV more accurately estimates prediction error than 10-fold CV.

Keywords: Air Pollution, Exposure Model, Machine Learning, Ozone, Wildfire

1. Introduction

Ground-level ozone is toxic to humans, animals and plants and contributes significantly to climate change as the third most important greenhouse gas [1, 2, 3, 4, 5, 6, 7]. Short-term exposure is linked to increased mortality [3], decreased respiratory function, exacerbation of chronic obstructive pulmonary disease (COPD), bronchitis, emphysema and asthma [8, 9, 10, 11]. Long-term exposure has been linked with respiratory and cardiovascular mortality [12, 13], decreased lung function [14] and the progression of emphysema [15].

Wildfires contribute to the formation of ozone in the lower atmosphere (troposphere) by releasing volatile organic compounds (VOCs) and nitrogen oxides (NO_x), which react in the presence of sunlight to form ozone [16, 17]. Fires upwind of large population centers can expose millions or even tens of millions of people to ozone and other pollutants [18]. Climate change is expected to intensify wildfires, which will likely increase the prevalence of wildfire-related ozone exposure.[17, 19]

The health effects of wildfire-induced ozone exposure are poorly understood.

17 A study of hospital admissions in Port, Portugal, in 2005 while wildfires were
18 burning nearby, indicated ozone was significantly associated with cardiovascu-
19 lar disease admissions, but not with respiratory admissions [20]. This analysis,
20 however, did not control for weather or land-use covariates. During a bushfire
21 in southeastern Australia, respiratory emergency department visits were signifi-
22 cantly associated with PM₁₀ (particulate matter 10 μm or smaller in diameter),
23 but not with ozone [21].

24 A key challenge facing epidemiological analyses of air pollution exposure is
25 quantifying pollution concentrations where people live, which may be distant
26 from regulatory monitoring sites. This is particularly difficult for wildfire-related
27 pollution, because wildfires often ignite far from urban regulatory monitoring
28 sites, and satellite evidence indicates that traditional monitoring networks are
29 too sparse to capture smoke plume variation and dynamics [22]. Epidemiologists
30 attempt to overcome this difficulty by constructing exposure models to predict
31 pollution concentration at unmonitored locations and times. The prediction of a
32 quantity across a domain, such as a spatial region, based on observations of that
33 quantity at discrete locations within that domain is referred to as downscaling
34 or interpolation, and is an example of infill prediction.

35 The simplest downscaling exposure models rely upon the tendency of nearby
36 observations to be more similar than those farther apart to interpolate between
37 pollution monitor observations without the use of additional information (i.e.,
38 without covariates). This tendency is an example of spatial (or space-time)
39 dependence. Kriging is a very commonly used method for interpolating ob-
40 servations, modeling air pollution concentrations as the best linear unbiased
41 prediction (BLUP) given the data and mean and covariance functions selected
42 by the researcher and estimated from the data [23]. Kriging tends to perform
43 relatively well when monitoring data is dense, but model accuracy degrades at
44 locations or times distant from monitor observations.

45 To improve accuracy, especially when monitoring density is sparse, researchers
46 have employed regression models that incorporate ancillary information as co-
47 variates. These models are referred to as land use regression in the literature,

48 but the covariates need not pertain to land use, and increasingly include satel-
49 lite retrievals, meteorological data, and less frequently the output of atmospheric
50 chemistry numerical simulation models. Land use regression models have been
51 used for modeling air pollution exposure at least since the Small Area Variations
52 in Air quality and Health (SAVIAH) study in 1997 [24] with numerous exam-
53 ples appearing subsequently. While these regression models include covariate
54 information, they have often assumed linear, additive covariate effects. This
55 assumption makes the effects easy to interpret but is too stringent to predict air
56 pollution concentrations accurately. These models cannot accommodate non-
57 linear effects or interactions between covariates unless they are specified by the
58 analyst a priori. They also lack a mechanism for variable selection, requiring
59 the analyst to manually select covariates or employ a separate variable selection
60 procedure.

61 These limitations have prompted the development of more flexible models
62 that allow for nonlinear effects including spatially or spatiotemporally varying-
63 coefficient models [25]. These models generally fit covariate effects with smooth
64 functionals that need not be linear and may be indexed by space or space-time,
65 which allows the covariate effects to differ across space and time. These models
66 are an improvement over the very stringent restrictions set on covariate effects
67 in linear, additive models, but they often rely on research code, making them
68 less accessible to other researchers. In our experience we have also found that
69 these more sophisticated models do not scale well to realistic space-time data
70 settings.

71 These challenges have motivated researchers to turn to more flexible models
72 that do not require such stringent assumptions, including a variety of machine
73 learning algorithms, which have been shown to be very useful for prediction [26],
74 especially random forest [27], gradient boosting [28] and neural networks [29].
75 They may lack the straightforward interpretability of linear regression, but this
76 is of secondary concern when prediction is the primary objective, and variable
77 importance scores have been developed for many such methods to quantify the
78 contribution of each covariate.

79 Generalized additive models [30], support vector machines [31, 32], gradi-
80 ent boosting [33] and deletion/substitution/addition [34] have been used to
81 model particulate matter exposure. Neural networks [35, 36] and random forest
82 [37, 38, 39] have been used to model both particulate matter and ozone concen-
83 trations. A comparison of 11 machine learning models indicated that random
84 forest, gradient boosting and bagged trees predict $\text{PM}_{2.5}$ (particulate matter
85 smaller than $2.5 \mu\text{m}$ in diameter) concentrations well during a wildfire event
86 [22]. Random forest, boosting and Cubist performed well in a comparison of
87 8 machine learning tools predicting $\text{PM}_{2.5}$ in British Columbia [40]. Here we
88 conduct a similar analysis using ten machine learning algorithms to model ozone
89 exposure during a wildfire air pollution event for the first time, evaluating their
90 predictive accuracy for use as land use regression models.

91 Comparing and evaluating prediction models for dependent data is chal-
92 lenging. Cross-validation (CV) and the bootstrap are commonly used model
93 evaluation procedures that repeatedly fit a model to a training subset of the
94 data and evaluate the accuracy of its predictions on a different, test subset,
95 combining the performance across multiple test subsets into a nonparametric
96 estimate of prediction error. For data that are spatially and temporally depen-
97 dent (i.e., autocorrelated), however, these procedures can be overly optimistic,
98 because of the dependence between training and test subsets [41].

99 In the case of daily air pollution monitoring observations, we wish to es-
100 timate the average error made by a downscaling model when predicting at a
101 new location within the spatial domain of the data. Including observations in
102 the training data that were taken at the same monitors as the test set obser-
103 vations provides an unrealistic amount of information on the test data, because
104 of the strong correlation between observations taken at the same location. This
105 produces estimates of prediction error that are biased downward, especially for
106 flexible models which tend to overfit to a greater degree than less flexible mod-
107 els when trained on dependent data. The true prediction error associated with
108 predicting at a new location would be greater than these estimates, because
109 the model could not have been trained on any observations recorded at a new

110 location.

111 When dependence is restricted to observations within the same group or
112 cluster, consistent (i.e., asymptotically unbiased) estimates of prediction error
113 can be recovered by resampling groups rather than individual observations. It is
114 unrealistic, however, to assume that spatially dependent data are nested within
115 independent groups of observations. Modified cross-validation schemes have
116 been used on pollution exposure data that partition the data into spatial grid
117 cells [33] or monitor locations [42, 43]. Such approaches attempt to reduce the
118 dependence between training and test data sets by placing all observations at a
119 particular location or within a particular region into the same cross-validation
120 fold. In this vein, we use leave-one-location-out (LOLO) cross-validation, which
121 defines each CV fold as the observations recorded at a single monitor location
122 [44]. This does not partition the data into independent groups, but estimates
123 the error associated with predicting the time series of ozone observations at a
124 new location, conditioning upon the observed monitor data. By using all the
125 observations at a single location as the test set, LOLO CV ensures that no ob-
126 servations from this location appear in the training set, which would result in
127 unrealistically low estimates of prediction error. It also avoids the overly conser-
128 vative bias of leave- k -locations-out CV, which tends to overestimate prediction
129 error because it uses substantially fewer observations for model training.

130 **2. Materials and Methods**

131 *2.1. Data*

132 One hundred ground-based ozone monitors administered by the United States
133 Environmental Protection Agency (EPA) made hourly observations from which
134 the daily maximum 8-hour averages were computed across northern California
135 between May 6, 2008 and September 26, 2008 for a total of 13,487 observa-
136 tions. We selected this time period with the goal of estimating ozone exposures
137 before, during, and after a spate of wildfires that afflicted northern California
138 in late June and July of 2008. The mean maximum 8-hour average concen-

Table 1: Covariates used to predict ozone.

Covariate	Data Source
Monitor Latitude	U. S. Environmental Protection Agency
Monitor Longitude	U. S. Environmental Protection Agency
Elevation (m)	National Digital Elevation Model
Date	U. S. Environmental Protection Agency
Dew Point ($^{\circ}$ K)	Rapid Update Cycle
Boundary Layer Height (m)	Rapid Update Cycle
Surface Pressure (Pa)	Rapid Update Cycle
Relative Humidity (%)	Rapid Update Cycle
Temperature at 2 m ($^{\circ}$ K)	Rapid Update Cycle
U-Component of Wind Speed (m/s)	Rapid Update Cycle
V-Component of Wind Speed (m/s)	Rapid Update Cycle
Inverse Distance to Nearest Fire (m^{-1})	Fire Inventory from NCAR v1.5
Annual Average Traffic within 1 km	Dynamap 2000, TeleAtlas
Agricultural Land Use within 1 km (%)	2006 National Land Cover Database
Urban Land Use within 1 km (%)	2006 National Land Cover Database
Vegetation Land Use within 1 km (%)	2006 National Land Cover Database
Normalized Difference Vegetation Index	Landsat Data
Nitrogen Dioxide (\log molecules/ cm^2)	Ozone Monitoring Instrument Satellite
WRF-Chem Carbon Monoxide (\log moles/day)	WRF-Chem
WRF-Chem PM _{2.5} (\log kg/day)	WRF-Chem
WRF-Chem Ozone (\log 8 Hour Maximum)	WRF-Chem

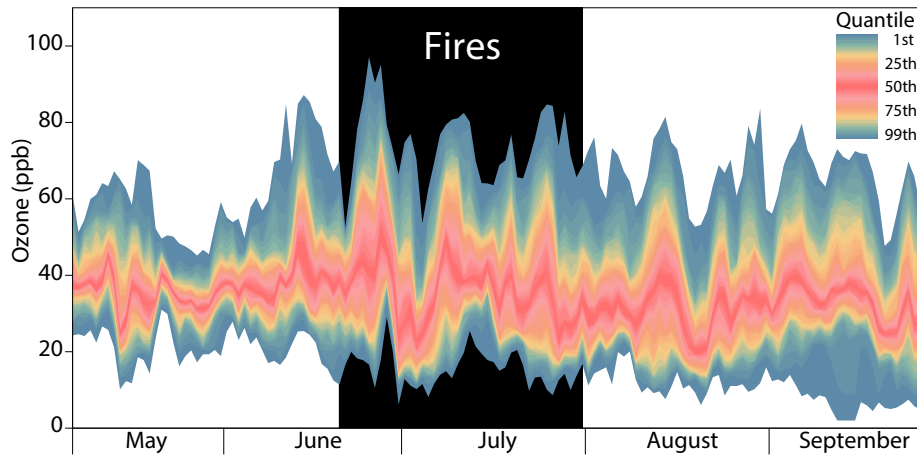


Figure 1: The daily empirical distribution of maximum daily 8-hour average ozone between May 6 and September 26, 2008 at 100 northern California monitoring sites.

139 tration was 36.2 ppb, and the standard deviation was 13.6 ppb. During the
 140 study, the maximum 8-hour average exceeded 70 ppb 236 times and exceeded
 141 75 ppb 107 times. Most exceedances occurred while the fires were burning (153
 142 and 75 respectively), although this time period is one of high solar intensity
 143 when high concentrations of ozone are expected. It is not our objective, how-
 144 ever, to quantify the contribution of wildfires to ozone formation, but simply to
 145 downscale ozone concentrations during a wildfire event for subsequent epidemio-
 146 logical analysis. Figure 1 depicts the temporal evolution of monitor observations
 147 throughout this time period.

148 Twenty-one covariates were also collected for the monitor locations, includ-
 149 ing location, elevation, date, atmospheric weather data (dew point, boundary
 150 layer height, surface pressure, relative humidity, temperature, and wind speed),
 151 inverse distance to the nearest fire, traffic, land use information (agricultural,
 152 urban, and vegetation), tropospheric nitrogen dioxide (NO_2) vertical column
 153 density and predictions of daily total carbon monoxide concentration (CO),
 154 particulate matter ($\text{PM}_{2.5}$) and daily maximum 8-hour average ozone. Table 1
 155 lists the covariates and their sources.

156 Monitor elevation was determined from the 2010 National Elevation Dataset

157 for California. The date of each observation was encoded as the continuous co-
158 variate, Julian date. The U.S. National Centers for Environmental Prediction’s
159 Rapid Update Cycle atmospheric prediction model provided hourly predictions
160 of dew point, planetary boundary layer height, surface pressure, relative hu-
161 midity, temperature, and the U and V components of wind speed, which were
162 averaged into daily values [45].

163 Inverse distance to the nearest fire was included as a covariate. The Fire
164 Inventory from NCAR (FINN) v1.5 provided estimates of fire point locations
165 in California during the study period [46]. Fire points occurring within 5 km
166 of each other were clustered and circumscribed by a polygon using the ArcGIS
167 Aggregate Points tool, and the distance between each monitoring site and the
168 closest point on the nearest fire cluster polygon was determined on each day
169 using the ArcGIS Near tool. On days with no fire in California, distance to the
170 nearest fire was undefined. Conceptualizing this undefined distance as equivalent
171 to the nearest fire cluster being infinitely far away, inverse distance to fire was
172 defined as 0 for observations taken on days with no fires in California and as
173 the inverse of the distance to the nearest fire cluster otherwise.

174 Dynamap 2000, a TeleAtlas product, was used to compute the annual av-
175 erage of roadway traffic within 1 km of each monitor [22]. The National Land
176 Cover Database for 2006 [47] was used to calculate the percentage of urban
177 development (codes 22, 23, and 24), agriculture (codes 81 and 82) and other
178 vegetation (codes 21, 41, 42, 43, 52, and 71) within 1 km of each monitor.

179 The normalized difference vegetation index (NDVI) quantifies the density
180 of green vegetation on a scale between -1 and 1 by measuring the visible and
181 near-infrared light reflected at a location via remote sensing. The chlorophyll
182 in healthy vegetation absorbs most of the visible light and reflects much of the
183 near-infrared light to which it is exposed, giving locations with more vegetation
184 a higher NDVI score. NDVI for each monitor location was extracted from the
185 NDVI remote sensing raster surface and included as a covariate.

186 Nitrogen dioxide (NO_2) was estimated on each day at monitor locations
187 (if available) using the Berkeley High-Resolution (BEHR) NO_2 tropospheric

188 column density retrieved from NASA’s Ozone Monitoring Instrument (OMI)
189 satellite, which has an overpass time of 1:30 local time [48] and a resolution
190 varying between 13 x 24 km to 42 x 162 km.

191 Predictions of daily total carbon monoxide (CO) and PM_{2.5} and the maxi-
192 mum daily 8-hour average ozone concentration were extracted for each day from
193 the Weather Research and Forecasting with Chemistry (WRF-Chem) 3.2 model.
194 WRF-Chem is a regional chemical transport model that simulates meteorology
195 and behavior of atmospheric gases and aerosols [49, 50]. Appendix A in the
196 supplemental material details the WRF-Chem inputs and options used for our
197 simulations.

198 2.2. Statistical Analysis

199 Each observation comprises an outcome, y_i , the log maximum 8-hour average
200 ozone on a given day at a given monitoring location, and a vector of covariates,
201 $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})'$, $i = 1, \dots, n$, where n is the number of observations, and p is
202 the number of covariates. The vector of outcomes, $\mathbf{y} = (y_1, \dots, y_n)'$, and the
203 matrix of covariates, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$, together compose the data, $D = \{\mathbf{X}, \mathbf{y}\}$.
204 Ozone observations were log transformed to reduce the impact of heteroscedas-
205 ticity (non-constant variance across the range of a variable), as data exploration
206 revealed the variance was substantially greater than the mean at high values.
207 The maximum daily 8-hour average ozone from the WRF chemical transport
208 model (WRF-Chem) was also log transformed to have the same scale as the out-
209 come. All other covariates were transformed to have a mean of 0 and variance
210 of 1.

211 Ten predictive algorithms were trained and evaluated on these data: elastic
212 net regression, generalized additive models (GAM), gradient boosting, k -nearest
213 neighbor regression, lasso regression, linear models, multivariate adaptive re-
214 gression splines (MARS), neural network, random forest, and support vector
215 machines with a radial basis kernel (SVM). All of the models except for neural
216 networks were fit using models available in version 6.0 of the caret R package
217 [51]. Neural networks were given special consideration on account of their grow-

218 ing popularity as machine learning prediction tools and especially the recent
 219 publication of papers using a neural network with inverse distance weighted con-
 220 volutional layers to predict ozone and particulate matter [35, 36]. We tested a
 221 neural network that mimicked those models, employing inverse distance weight-
 222 ing to create convolutional spatial, temporal and space-time layers using the
 223 keras R package [52]. These models were less accurate than a standard feedfor-
 224 ward neural network, and so we have reported the results of that network here.
 225 Training each prediction model produces a prediction rule $\eta(\mathbf{x}, D_T)$, which is a
 226 function of D_T , the data on which it was trained, and a vector of covariates, \mathbf{x} ,
 227 mapping them to a prediction for $y \mid \mathbf{x}$, which is often used as an estimator of
 228 $E(y \mid \mathbf{x})$, the conditional expectation of y given \mathbf{x} .

229 The models were tuned, selected, and evaluated using cross-validated esti-
 230 mators of root mean square error (RMSE) and R^2 , which are both functions of
 231 the mean square error (MSE). The MSE of a prediction rule $\eta(\mathbf{x}, D_T)$, where
 232 D_T is the data with which η was trained, may be estimated using a test data
 233 set D_W as

$$\hat{MSE}(D_W, \eta(\mathbf{x}, D_T)) = \frac{1}{n_w} \sum_{j \in D_W} (y_j - \eta(\mathbf{x}_j, D_T))^2, \quad (1)$$

234 where n_w is the number of data points in D_W . If D_W and D_T are disjoint,
 235 (i.e., if η was not trained using any part of D_W), then this is an out-of-sample
 236 estimator of the MSE. RMSE may be estimated by the square root of \hat{MSE} ,
 237 and R^2 is estimated by

$$\hat{R}^2(D_W, \eta(\mathbf{x}, D_T)) = 1 - \frac{\hat{MSE}(D_W, \eta(\mathbf{x}, D_T))}{n_w^{-1} \sum_{j \in D_W} (y_j - \bar{y}_w)^2}, \quad (2)$$

238 where $\bar{y}_w = n_w^{-1} \sum_{j \in D_W} y_j$ is the mean outcome in D_W . For ease of notation,
 239 the function arguments for \hat{MSE} , \hat{RMSE} , and \hat{R}^2 are hereafter suppressed.

240 Two different cross-validation (CV) strategies were employed for model eval-
 241 uation: 10-fold cross-validation and leave-one-location-out (LOLO) cross-validation.
 242 For 10-fold CV, the data were randomly partitioned into 10 non-overlapping
 243 subsets, each containing one tenth of the data. Each subset served as the test

244 data for models trained on the other nine tenths of the data, resulting in ten
245 different pairs of training and test sets, with each observation appearing in one
246 test set and the nine training sets not paired with that test set. This yielded 10
247 estimates of MSE for each model, which were averaged into an overall estimate
248 of MSE, from which the 10-fold CV estimates of RMSE and R^2 were computed.

249 Ten-fold cross-validation is widely used for estimating prediction error; how-
250 ever, it is known to be overly optimistic for dependent data [41]. Data recorded
251 by air pollution monitors are expected to exhibit spatial or space-time depen-
252 dence. To more accurately estimate the downscaling error associated with pre-
253 dicting ozone at an unobserved location, RMSE and R^2 were estimated using
254 LOLO CV, in which a model is trained on data from all but one location, and its
255 prediction error is computed for the observations at the withheld location. This
256 process is repeated with observations at each location serving as the withheld
257 test set once, and the resulting errors are averaged into the LOLO CV estimate
258 of prediction error. Unlike 10-fold CV in which observations are distributed
259 among folds uniformly at random, LOLO CV ensures that no observations from
260 the test location may appear in the training data. This provides a realistic
261 estimate of the downscaling prediction error associated with predicting ozone
262 observations at a new location within the same region as the monitoring data.

263 Most predictive machine learning algorithms depend upon one or more pa-
264 rameters whose values must be set prior to fitting the model. Algorithm per-
265 formance can vary greatly depending on these parameter values, and it is often
266 desirable to select values that optimize some criteria in an attempt to improve
267 model performance. The process of choosing values for these parameters is often
268 referred to as tuning and the parameters themselves as tuning parameters (or
269 hyperparameters). In our analysis, most tuning parameter values were selected
270 by comparing the performance of candidate values on 25 bootstrap samples of
271 the data using the caret R package [51]. Parameters for k -nearest neighbors and
272 GAM were specifically tuned for LOLO CV in an attempt to stabilize the LOLO
273 CV prediction error, as these models made extremely poor LOLO predictions
274 using bootstrap-selected tuning parameter values. Appendix C in the supple-

275 mental material details these tuning procedures and their results. Substantial
276 effort was taken in selecting the number of layers, nodes, activation functions
277 and distance weighting functions for the neural network. The most accurate
278 model was a feedforward neural network with one hidden layer of 21 nodes
279 using a rectified linear unit activation function without the inverse distance
280 weighting functions and convolutional layers employed in previously published
281 models [35, 36].

282 To investigate the transferability of a model trained on data in one region
283 to another, i.e., its ability to extrapolate rather than downscale, the predictive
284 performance of the two best models trained on the 2008 northern California
285 wildfire period—those two with the lowest LOLO CV estimates of RMSE—was
286 evaluated on data collected during a 2007 wildfire event in southern California.
287 The southern California data consisted of 5,978 daily 8-hour maximum ozone
288 values recorded at 72 monitors between September 1, 2007 and November 28,
289 2007.

290 3. Results and Discussion

291 Figure 2 graphically depicts the cross-validated estimates of RMSE and R^2
292 for each algorithm using 10-fold CV and LOLO CV. In every case, the 10-fold
293 CV $RM\hat{S}E$ was lower than the LOLO CV $RM\hat{S}E$, and the 10-fold CV \hat{R}^2 was
294 higher than the LOLO CV \hat{R}^2 . Gradient boosting had the lowest 10-fold CV
295 $RM\hat{S}E$ (0.186 log ppm), lowest LOLO CV $RM\hat{S}E$ (0.228 log ppm), highest
296 10-fold CV \hat{R}^2 (0.784), and highest LOLO CV \hat{R}^2 (0.677). Random forest
297 placed second in all four categories. The table in Appendix B lists exact values
298 for $RM\hat{S}E$ and \hat{R}^2 for each model. These results answer the two primary
299 questions posed by this study, demonstrating that machine learning methods
300 can downscale ozone during a wildfire with reasonable accuracy and identifying
301 gradient boosting and random forest as performing particularly well.

302 The 10-fold CV estimates of RMSE and R^2 were optimistic compared to
303 those of LOLO CV for all ten models. This over optimism of 10-fold CV is intu-

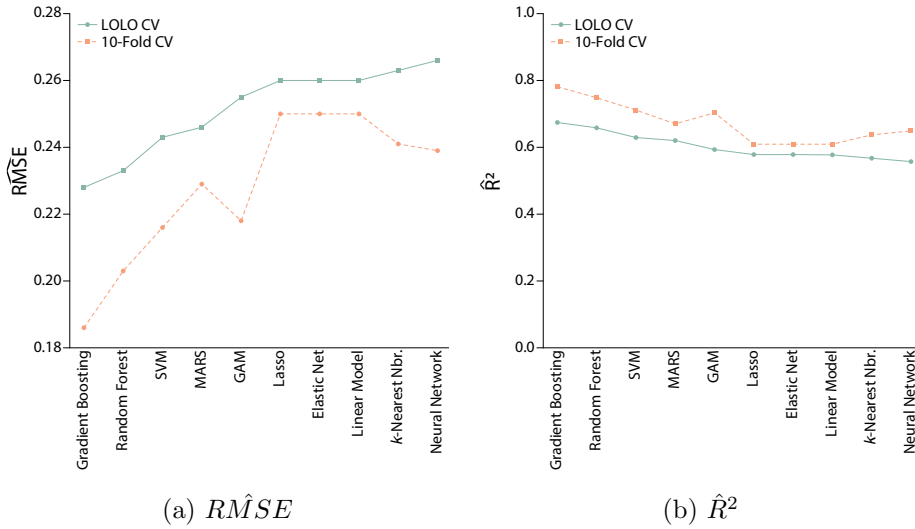


Figure 2: 10-fold and leave-one-location-out cross-validated estimates of RMSE and R^2 for downscaling ozone prediction models.

304 itive, because of the strong dependence between test and training observations
 305 from the same monitor location. The 10-fold CV estimators are also unreliable
 306 for model selection. The ordering of model performance is not invariant to the
 307 choice of evaluation criterion, which is demonstrated here by the evaluation of
 308 the neural network. It is the worst performing model when evaluated by LOLO
 309 CV, but is sixth best according to 10-fold CV. The ordering of GAM, MARS
 310 and k -nearest neighbors also differ, though the magnitude of those differences
 311 is not as substantial.

312 The difference between the 10-fold and LOLO cross-validated estimates of
 313 performance was smaller for relatively inflexible models like lasso, elastic net,
 314 and linear regression than for the other models, whose greater flexibility enabled
 315 them to better exploit the more highly dependent folds of 10-fold CV. The large
 316 difference for k -nearest neighbor regression is due to the strong dependence
 317 between observations recorded at the same monitor location. In 10-fold CV,
 318 the nearest neighbors of an observation are very likely to be other observations
 319 taken at that location. In LOLO CV, no observations taken at the test location

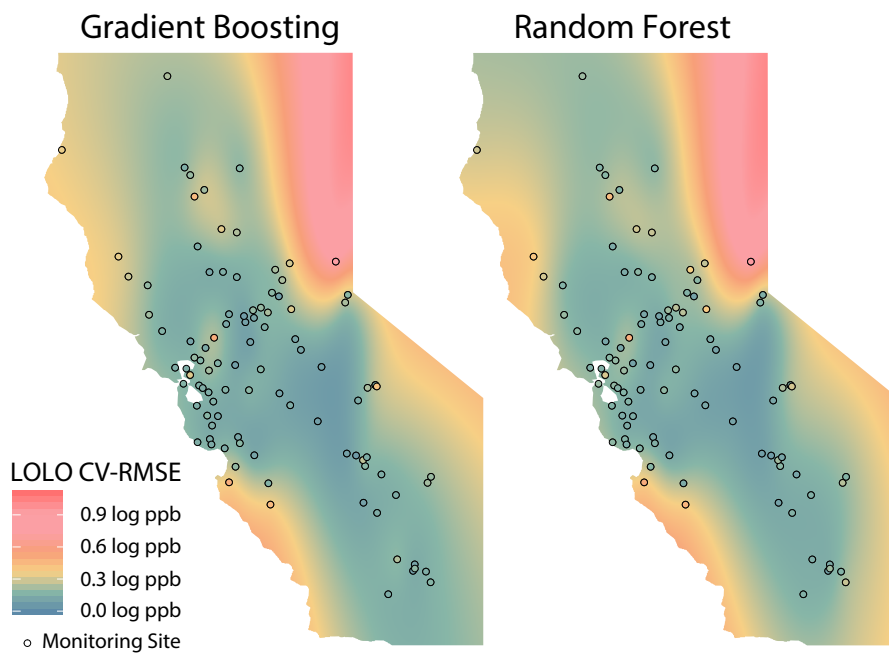


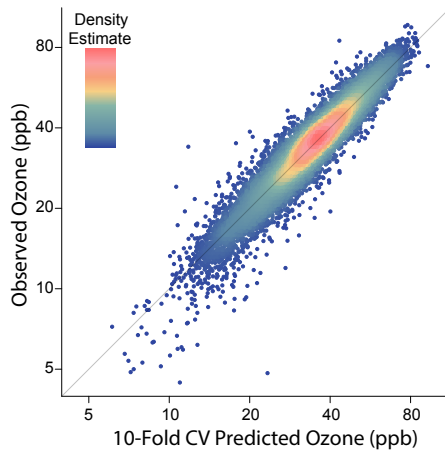
Figure 3: The leave-one-location-out (LOLO) cross-validated estimates of RMSE averaged over the study period (May 6, 2008–September 26, 2008) are plotted at each monitor location for gradient boosting (left) and random forest (right). These average prediction error estimates were then smoothed throughout the study region using a two-dimensional spline-on-sphere smoother to provide a visual estimate of how downscaling predictive performance may vary across the spatial domain.

320 appear in the training data. It is no surprise that this yields substantially higher
321 estimates of prediction error.

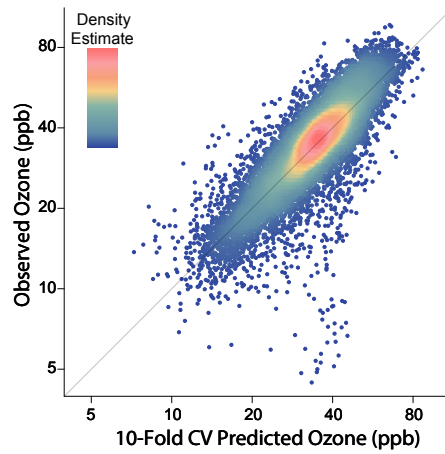
322 Figure 4 plots 10-fold and LOLO CV predictions against observed daily 8-
323 hour maximum average ozone for gradient boosting and random forest. The CV
324 prediction for each point is the predicted value when it appears in the test fold
325 of the CV procedure. Points that fall on the grey diagonal line are perfectly
326 predicted. The tighter clustering of points around this line in the 10-fold CV
327 plots corresponds to the more accurate predictions made when test set monitor
328 locations are included in the training data.

329 The neural network performed less well than in previous applications for
330 predicting ozone and particulate matter over a spatial grid across the continental
331 United States [35, 36]. The recent popularity of convolutional neural networks
332 is largely due to their performance on image-processing problems. Gridded
333 spatial (or space-time) data bear a much greater resemblance to image data
334 than do the point process monitor data upon which they were evaluated here.
335 It is also possible that alternate specifications of the network architecture could
336 improve performance, but developing such a model goes well beyond the scope
337 of this comparison, which is limited to readily available algorithms that do not
338 require substantial expertise in model specification or implementation. At least
339 in this context neural networks do not succeed as an automatically regularized
340 statistical learning tool (i.e., as a black box), but their performance in other
341 studies suggests they may work well as a highly specialized tool designed using
342 domain knowledge specific to a particular application.

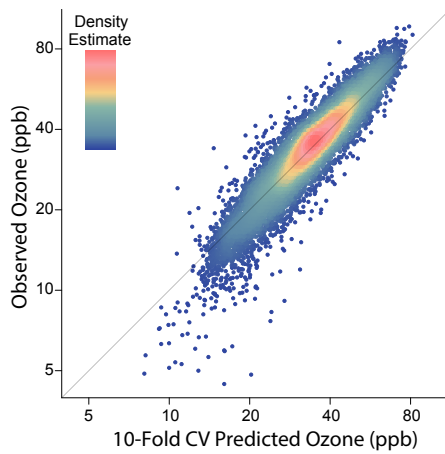
343 The magnitude of the difference between the LOLO and 10-fold CV esti-
344 mates of prediction error has meaningful consequences for estimating exposures
345 for subsequent epidemiological analyses. Downscaled exposure is often used as
346 the covariate of interest in analyses seeking to infer the health consequences
347 of air pollution without accounting for prediction uncertainty. The more real-
348 istic estimates of prediction error provided by LOLO CV offer better insight
349 into whether it is reasonable to ignore this uncertainty. This may motivate
350 improvements to epidemiological models to account for exposure measurement



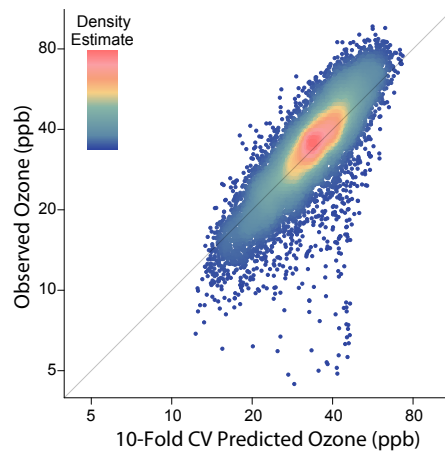
(a) Gradient Boosting 10-Fold CV



(b) Gradient Boosting LOLO CV



(c) Random Forest 10-Fold CV



(d) Random Forest LOLO CV

Figure 4: 10-fold and leave-one-location-out (LOLO) cross-validated gradient boosting and random forest predictions plotted against observed daily 8-hour maximum average ozone on the log scale.

351 error.

352 The increased accuracy of LOLO CV comes at a computational cost. When
353 the number of locations exceeds 10, LOLO CV is more computationally ex-
354 pensive than 10-fold CV. In this analysis, there are 100 monitor locations and
355 therefore 100 folds in LOLO CV, corresponding to approximately 10 times the
356 computational burden of 10-fold CV. Grouping monitor locations into folds is
357 an appealing strategy to alleviate this burden [43], however, it estimates pre-
358 diction error over a different spatial resolution than LOLO CV, and will result
359 in overly conservative estimates.

360 The top two models, gradient boosting and random forest, are both ensem-
361 bles of tree-based models that provide very flexible mean structures. Their excel-
362 lence suggests that the mean structure characterizing the relationship between
363 covariates and ozone likely includes interactions, non-linearities and possibly
364 discontinuities. These results do not demonstrate that the underlying chemi-
365 cal processes by which ozone forms are similarly complicated, but that seems
366 likely. The 10-fold CV estimates of RMSE and R^2 were similar to, although
367 slightly lower than, those reported in a similar analysis of machine learning ex-
368 posure models for $PM_{2.5}$ during the same wildfire time period [22]. Tree-based
369 ensembles were also the best performing models in that study, suggesting that
370 algorithms with flexible mean structures can produce useful exposure models
371 for ozone and $PM_{2.5}$ during wildfire events. Traditional exposure models have
372 focused on modeling the dependence between observations, while employing a
373 simple mean structure. The machine learning models evaluated here assume in-
374 dependent observations, but offer much greater flexibility in modeling the mean.
375 This approach is expected to provide more accurate predictions distant from
376 the observations on which the model was trained than methods that rely upon
377 the dependence between observations. Combining the flexible mean structure
378 of tree-based ensembles with the dependence structures of traditional spatial
379 statistics models is a promising avenue for future work.

380 An interesting related analysis would be assessing the effect wildfires have
381 on ozone formation. Doing so would require additional numerical simulations

382 from the WRF-Chem model that exclude wildfire emissions from its inputs.
383 With the data currently available to us, it is impossible to disentangle the effect
384 of wildfires from the other inputs of the WRF-Chem model. This inferential
385 problem is also quite different from the downscaling task which is the focus of
386 this study and may require an entirely different modeling approach. A model
387 that provides excellent downscaling predictions may not be useful for drawing
388 scientific inference.

389 Another interesting question is whether ozone formation was NO_x-limited
390 or VOC-limited. During a wildfire the chemical regime is primarily determined
391 by the amount and conditions of the wildfire fuel, leading to rapid changes in
392 NO_x and VOC sensitivity from day to day and even within the course of a
393 day. Understanding this would require a fully separate chemical analysis of air
394 quality conditions in northern California that is beyond both the scope of this
395 paper and the scope of our data, as we lack data on VOC concentrations. One
396 previous study, however, examined the changes that occurred in atmospheric
397 chemistry when wildfire plumes interacted with urban pollution during these
398 fires [53].

399 We also lack information on chlorofluorocarbon (CFC) emissions, which
400 break down ozone and thus influence ozone concentrations. If data on these
401 compounds were available during the study period, we could include them as
402 covariates in an attempt to improve predictions. The absence of data on VOC
403 and CFC emissions does not invalidate the downscaling enterprise. Downscal-
404 ing models are constructed to combine the available information, whether from
405 observational processes or complex computational models like WRF-Chem, into
406 accurate predictions within a particular domain. They are not scientific models
407 and do not necessarily imply anything about the chemical or physical mecha-
408 nisms by which pollutants form and move. Using flexible models strictly for
409 downscaling also protects us from biases in the WRF-Chem output. We need
410 not validate the accuracy of the WRF-Chem simulations; we simply rely on
411 the models to learn the relationship between this output and observed ozone
412 concentrations.

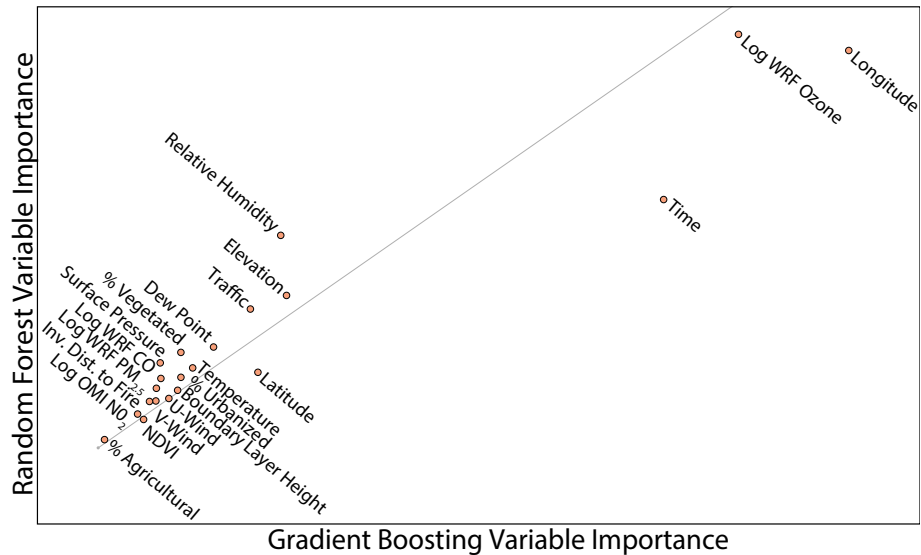


Figure 5: Pairwise normalized gradient boosting and random forest variable importance scores for models trained on the full data. The light grey line denotes equal importance in the two models.

413 Figure 5 plots covariate importance scores for random forest and gradient
 414 boosting models fit to the full data. Random forest variable importance was cal-
 415 culated as the mean decrease in residual sum of squares resulting from splitting
 416 on that covariate averaged across all the trees of the forest. Variable impor-
 417 tance for gradient boosting was calculated by permuting that covariate’s values
 418 and computing the average difference in MSE between predictions made with
 419 permuted and un-permuted values [54]. The variable importance scores for the
 420 two models were normalized to sum to one for ease of comparison. Most covari-
 421 ates are close to the grey diagonal, which indicates equal importance in the two
 422 models. Longitude, WRF-Chem ozone and time were the three most important
 423 covariates for both models. WRF-Chem is constructed to estimate atmospheric
 424 ozone and so it is not surprising that it has a high importance score. Longitude,
 425 latitude and time can proxy for unobserved factors, but also calibrate the effect
 426 of other covariates similar to space- or time-varying coefficient models. This
 427 calibration may be especially important for the numerical outputs of the WRF-

428 Chem model. Longitude is likely particularly important because it can be used
429 to index many of the significant geographical features of northern California
430 that run approximately North-South, including the coast, San Joaquin Valley,
431 coastal and Sierra Nevada mountain ranges. These geographical features may
432 be associated with important, unobserved information that is not captured by
433 the other covariates including VOC and CFC concentrations.

434 Neither model when trained on the northern California 2008 wildfire data
435 accurately predicted ozone exposure in southern California in 2007. The pre-
436 dictions from both models had negative \hat{R}^2 , indicating that their predictions
437 were less accurate (i.e., had higher estimated MSE) than the sample mean of
438 the southern California ozone monitors, which by definition has an \hat{R}^2 of 0. In
439 fairness to gradient boosting and random forest, in an out-of-domain prediction
440 problem, the sample mean is unknown, and therefore cannot be used as a pre-
441 diction rule. When downscaling gradient boosting and random forest models
442 were fit to the 2007 southern California wildfire data, they had LOLO CV errors
443 comparable to the LOLO CV errors reported above. These models accurately
444 downscaled ozone observations during the southern California wildfire (just as
445 they did for the northern California data), but they did not extrapolate well
446 outside of the domain in which they were trained.

447 This is not surprising and illustrates the proper interpretation of our mod-
448 eling efforts, which is statistical downscaling (i.e., interpolating) within the ob-
449 served space-time domain. The substantial decrease in predictive accuracy be-
450 tween within-domain (i.e., downscaling) and out-of-domain (i.e., extrapolating)
451 predictive performance suggests that the relationships between covariates and
452 ozone exposure differ in space and time and demonstrates the dangers of using a
453 downscaling model for extrapolation. Within the observed space-time domain,
454 space and time can proxy for unobserved spatially-indexed covariates in flex-
455 ible models like gradient boosting and random forest, improving downscaling
456 predictions, but impeding straightforward extrapolation to different space-time
457 domains where the relationship between these unmeasured covariates and space-
458 time may be different. One referee suggested that the difference in ozone pre-

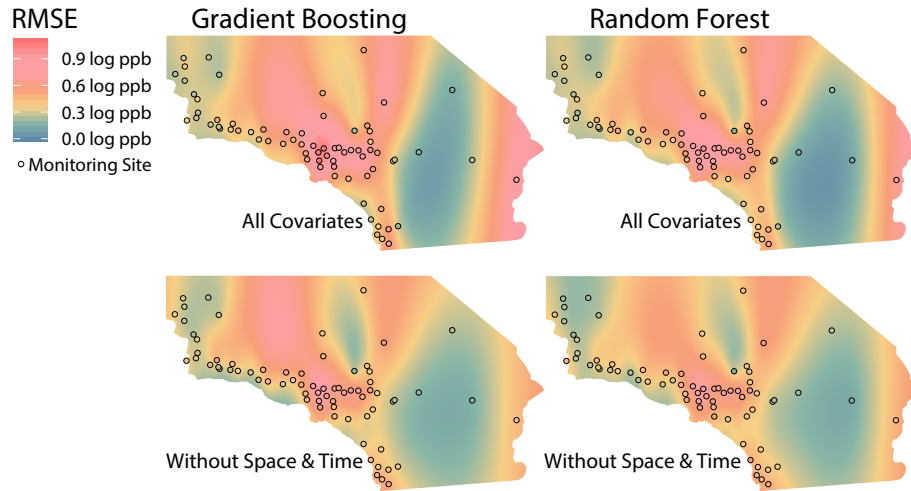


Figure 6: The extrapolation error made by gradient boosting and random forest models trained on the 2008 data and predicting ozone during the 2007 southern California fire. The mean RMSE at each location is smoothed throughout the study region using a two-dimensional spline-on-sphere smoother.

459 cursors between regions, specifically whether ozone formation is NO_x-limited or
 460 VOC-limited, is likely one such unobserved, spatially-indexed covariate hinder-
 461 ing spatial extrapolation.

462 As a further check, we repeated this extrapolating procedure excluding lati-
 463 tude, longitude and time from the covariates. The domains of the other covari-
 464 ates were comparable between the two data sets, with the exception of a few
 465 low values for surface pressure in the southern California data. Excluding spa-
 466 tial and temporal covariates improved predictive accuracy, reducing LOLO CV
 467 \hat{RMSE} from 0.587 to 0.499 for gradient boosting and 0.544 to 0.462 for random
 468 forest. As expected, space-time covariates improve downscaling predictions but
 469 worsen extrapolating predictions, because the unobserved information indexed
 470 by these covariates differs in other regions and at different times.

471 The scale on which prediction is performed is also important and application
 472 specific. In this analysis we performed prediction on the log scale to normal-
 473 ize the variance, but also because it balances large and small errors allowing

474 accurate predictions to be made at both large and small concentrations. We
475 believe this is a natural scale for the intended subsequent applications of our
476 predictions in epidemiological analyses of pollution health effects. Prediction
477 on the untransformed, original scale would more heavily weight predictions at
478 high ozone concentrations. This may be useful for some applications, but we
479 prefer to balance predictive accuracy at low and high concentrations for our
480 application.

481 The comparison performed here demonstrates that machine learning predic-
482 tion algorithms, especially ensembles of tree models like gradient boosting and
483 random forest, can accurately downscale ozone concentrations during wildfire
484 events. We believe they would downscale ozone similarly well in the absence
485 of wildfire events. The models we consider here, however, did not accurately
486 extrapolate beyond the space-time domain on which they were trained. This
487 analysis also demonstrates that the choice of evaluation metric is critical to un-
488 derstanding predictive performance. Metrics that ignore the dependent struc-
489 ture of the data, including k -fold CV, are overly optimistic and unreliable for
490 model selection. LOLO CV is a superior alternative that accounts of the spatial
491 dependence of the data in evaluating model predictive performance, resulting
492 in more reliable estimates of predictive performance.

493 **4. Acknowledgments**

494 This work was funded in part by a grant from the Bureau of Land Manage-
495 ment [grant number L14AC00173].

496 **References**

- 497 [1] World Health Organization, Health effects of particulate matter. Policy
498 implications for countries in eastern Europe, Caucasus and central Asia.
499 Copenhagen: WHO Regional Office for Europe; 2013.
- 500 [2] D. B. Menzel, Ozone: an overview of its toxicity in man and animals, J.
501 Toxicol. Environ. Health.

- 502 [3] M. L. Bell, A. McDermott, S. L. Zeger, J. M. Samet, F. Dominici, Ozone
503 and short-term mortality in 95 us urban communities, 1987-2000, JAMA
504 292 (19) (2004) 2,372–2,378.
- 505 [4] T. Stocker, D. Qin, G. Plattner, M. Tignor, S. Allen, J. Boschung,
506 A. Nauels, Y. Xia, V. Bex, P. Midgley, Fifth assessment report of the
507 intergovernmental panel on climate change, The Phys. Sci. Basis.
- 508 [5] K. R. Smith, M. Jerrett, H. R. Anderson, R. T. Burnett, V. Stone, R. Der-
509 went, R. W. Atkinson, A. Cohen, S. B. Shonkoff, D. Krewski, C. A.
510 Pope, Public health benefits of strategies to reduce greenhouse-gas emis-
511 sions: health implications of short-lived greenhouse pollutants, The Lancet
512 374 (9707) (2009) 2091–2103.
- 513 [6] R. A. Silva, J. J. West, J.-F. Lamarque, D. T. Shindell, W. J. Collins,
514 G. Faluvegi, G. A. Folberth, L. W. Horowitz, T. Nagashima, V. Naik, S. T.
515 Rumbold, Future global mortality from changes in air pollution attributable
516 to climate change, Nat. Clim. Change 7 (9) (2017) 647.
- 517 [7] D. Fowler, M. Amann, F. Anderson, M. Ashmore, P. Cox, M. Depledge,
518 D. Derwent, P. Grennfelt, N. Hewitt, O. Hov, M. Jenkin, Ground-level
519 ozone in the 21st century: future trends, impacts and policy implications,
520 R. Soc. Sci. Policy Rep. 15 (08).
- 521 [8] National Research Council, Estimating mortality risk reduction and eco-
522 nomic benefits from controlling ozone air pollution, National Academies
523 Press, 2008.
- 524 [9] S. Sousa, M. Alvim-Ferraz, F. Martins, Health effects of ozone focusing on
525 childhood asthma: what is now known—a review from an epidemiological
526 point of view, Chemosphere 90 (7) (2013) 2051–2058.
- 527 [10] M. Guarneri, J. R. Balmes, Outdoor air pollution and asthma, The Lancet
528 383 (9928) (2014) 1581–1592.

- 529 [11] J. Y. Lee, S.-B. Lee, G.-N. Bae, A review of the association between air
530 pollutant exposure and allergic diseases in children, *Atmos. Pollut. Res.*
531 5 (4) (2014) 616–629.
- 532 [12] M. C. Turner, M. Jerrett, C. A. Pope III, D. Krewski, S. M. Gapstur, W. R.
533 Diver, B. S. Beckerman, J. D. Marshall, J. Su, D. L. Crouse, R. T. Burnett,
534 Long-term ozone exposure and mortality in a large prospective study, *Am.*
535 *J. Respir. Crit. Care Med.* 193 (10) (2016) 1,134–1,142.
- 536 [13] D. L. Crouse, P. A. Peters, P. Hystad, J. R. Brook, A. van Donkelaar,
537 R. V. Martin, P. J. Villeneuve, M. Jerrett, M. S. Goldberg, C. A. Pope III,
538 M. Brauer, R. D. Brook, A. Robichaud, R. Menard, R. T. Burnett, Ambi-
539 ent pm_{2.5}, o₃, and no₂ exposures and associations with mortality over 16
540 years of follow-up in the canadian census health and environment cohort
541 (canche), *Environ. Health Perspect.* 123 (11) (2015) 1180.
- 542 [14] B.-F. Hwang, Y.-H. Chen, Y.-T. Lin, X.-T. Wu, Y. L. Lee, Relationship
543 between exposure to fine particulates and ozone and reduced lung function
544 in children, *Environ. Res.* 137 (2015) 382–390.
- 545 [15] R. Barr, E. Hoffman, J. Madrigano, C. Aaron, P. Sampson, L. Sheppard,
546 S. Vedal, J. Kaufman, M. Wang, Long-term exposure to ozone and acceler-
547 ated progression of percent emphysema and decline in lung function: The
548 mesa air and lung studies, *Medicine* 1 (2) (2016) 3.
- 549 [16] M. Val Martin, R. Honrath, R. C. Owen, G. Pfister, P. Fialho, F. Barata,
550 Significant enhancements of nitrogen oxides, black carbon, and ozone in the
551 north atlantic lower free troposphere resulting from north american boreal
552 wildfires, *J. Geophys. Res. D: Atmos.* 111 (D23).
- 553 [17] D. A. Jaffe, N. L. Wigder, Ozone production from wildfires: A critical
554 review, *Atmos. Environ.* 51 (2012) 1–10.
- 555 [18] J. Wu, A. M. Winer, R. J. Delfino, Exposure assessment of particulate

- 556 matter air pollution before, during, and after the 2003 southern california
557 wildfires, *Atmos. Environ.* 40 (18) (2006) 3333–3348.
- 558 [19] U. Confalonieri, B. Menne, R. Akhtar, K. L. Ebi, M. Hauengue, R. S.
559 Kovats, B. Revich, A. Woodward, Human health (2007), in: M. Parry,
560 O. Canziani, J. Palutikof (Eds.), *Climate change 2007: impacts, adap-
561 tation and vulnerability*, Vol. 4, Cambridge University Press Cambridge,
562 Cambridge, 2007.
- 563 [20] J. M. Azevedo, F. L. Gonçalves, M. de Fátima Andrade, Long-range ozone
564 transport and its impact on respiratory and cardiovascular health in the
565 north of portugal, *Int. J. Biometeorol.* 55 (2) (2011) 187–202.
- 566 [21] R. Tham, B. Erbas, M. Akram, M. Dennekamp, M. J. Abramson, The
567 impact of smoke on respiratory hospital outcomes during the 2002–2003
568 bushfire season, victoria, australia, *Respirology* 14 (1) (2009) 69–75.
- 569 [22] C. E. Reid, M. Jerrett, M. L. Petersen, G. G. Pfister, P. E. Morefield,
570 I. B. Tager, S. M. Raffuse, J. R. Balmes, Spatiotemporal prediction of
571 fine particulate matter during the 2008 northern california wildfires using
572 machine learning, *Environ. Sci. Technol.* 49 (6) (2015) 3887–3896.
- 573 [23] M. L. Stein, *Interpolation of spatial data: some theory for kriging*, Springer
574 Science & Business Media, 2012.
- 575 [24] D. J. Briggs, S. Collins, P. Elliott, P. Fischer, S. Kingham, E. Lebreton,
576 K. Pryn, H. Van Reeuwijk, K. Smallbone, A. Van Der Veen, Mapping urban
577 air pollution using gis: a regression-based approach, *Int. J. Geog. Inf. Sci.*
578 11 (7) (1997) 699–718.
- 579 [25] A. E. Gelfand, H.-J. Kim, C. Sirmans, S. Banerjee, Spatial modeling with
580 spatially varying coefficient processes, *J. Am. Stat. Assoc.* 98 (462) (2003)
581 387–396.
- 582 [26] C. Robert, *Machine learning, a probabilistic perspective*, Taylor & Francis,
583 2014.

- 584 [27] R. Caruana, A. Niculescu-Mizil, An empirical comparison of supervised
585 learning algorithms, in: Proceedings of the 23rd international conference
586 on Machine learning, ACM, 2006, pp. 161–168.
- 587 [28] A. J. Ferreira, M. A. Figueiredo, Boosting algorithms: A review of methods,
588 theory, and applications, in: Ensemble machine learning, Springer, 2012,
589 pp. 35–85.
- 590 [29] I. Goodfellow, Y. Bengio, A. Courville, Deep learning, MIT press, 2016.
- 591 [30] Y. Liu, C. J. Paciorek, P. Koutrakis, Estimating regional spatial and tem-
592 poral variability of $\text{pm}_{2.5}$ concentrations using satellite data, meteorology,
593 and land use information, *Environ. Health Perspect.* 117 (6) (2009) 886.
- 594 [31] W.-Z. Lu, W.-J. Wang, Potential assessment of the “support vector ma-
595 chine” method in forecasting ambient air pollutant trends, *Chemosphere*
596 59 (5) (2005) 693–701.
- 597 [32] H. Weizhen, L. Zhengqiang, Z. Yuhuan, X. Hua, Z. Ying, L. Kaitao,
598 L. Donghui, W. Peng, M. Yan, Using support vector regression to pre-
599 dict pm_{10} and $\text{pm}_{2.5}$, in: IOP Conference Series: Earth and Environmental
600 Science, Vol. 17, IOP Publishing, 2014, p. 012268.
- 601 [33] Y. Zhan, Y. Luo, X. Deng, H. Chen, M. L. Grieneisen, X. Shen, L. Zhu,
602 M. Zhang, Spatiotemporal prediction of continuous daily $\text{pm}_{2.5}$ concentra-
603 tions across china using a spatially explicit machine learning algorithm,
604 *Atmos. Environ.* 155 (2017) 129–139.
- 605 [34] B. S. Beckerman, M. Jerrett, M. Serre, R. V. Martin, S.-J. Lee,
606 A. Van Donkelaar, Z. Ross, J. Su, R. T. Burnett, A hybrid approach to esti-
607 mating national scale spatiotemporal variability of $\text{pm}_{2.5}$ in the contiguous
608 united states, *Environ. Sci. Technol.* 47 (13) (2013) 7,233–7,241.
- 609 [35] Q. Di, I. Kloog, P. Koutrakis, A. Lyapustin, Y. Wang, J. Schwartz, Assess-
610 ing $\text{pm}_{2.5}$ exposures with high spatiotemporal resolution across the conti-
611 nental united states, *Environ. Sci. Technol.* 50 (9) (2016) 4,712–4,721.

- 612 [36] Q. Di, S. Rowland, P. Koutrakis, J. Schwartz, A hybrid model for spatially
613 and temporally resolved ozone exposures in the continental united states,
614 *J. Air Waste Manage. Assoc.* 67 (1) (2017) 39–52.
- 615 [37] X. Hu, J. H. Belle, X. Meng, A. Wildani, L. A. Waller, M. J. Strickland,
616 Y. Liu, Estimating pm_{2.5} concentrations in the conterminous united states
617 using the random forest approach, *Environ. Sci. Technol.* 51 (12) (2017)
618 6,936–6,944.
- 619 [38] Y. Zhan, Y. Luo, X. Deng, M. L. Grieneisen, M. Zhang, B. Di, Spatiotem-
620 poral prediction of daily ambient ozone levels across china using random
621 forest for human exposure assessment, *Environmental Pollution* 233 (2018)
622 464–473.
- 623 [39] C. Brokamp, R. Jandarov, M. Hossain, P. Ryan, Predicting daily urban
624 fine particulate matter concentrations using a random forest model, *Envi-
625 ronmental science & technology* 52 (7) (2018) 4173–4179.
- 626 [40] Y. Xu, H. C. Ho, M. S. Wong, C. Deng, Y. Shi, T.-C. Chan, A. Knudby,
627 Evaluation of machine learning techniques with multiple remote sensing
628 datasets in estimating monthly concentrations of ground-level pm_{2.5}, *En-
629 vironmental Pollution* 242 (2018) 1417–1426.
- 630 [41] D. R. Roberts, V. Bahn, S. Ciuti, M. S. Boyce, J. Elith, G. Guillera-Arroita,
631 S. Hauenstein, J. J. Lahoz-Monfort, B. Schröder, W. Thuiller, D. I. Warton,
632 B. A. Wintle, F. Hartig, C. F. Dormann, Cross-validation strategies for data
633 with temporal, spatial, hierarchical, or phylogenetic structure, *Ecography*
634 40 (8) (2017) 913–929.
- 635 [42] M. Lee, I. Kloog, A. Chudnovsky, A. Lyapustin, Y. Wang, S. Melly,
636 B. Coull, P. Koutrakis, J. Schwartz, Spatiotemporal prediction of fine par-
637 ticulate matter using high-resolution satellite images in the southeastern
638 us 2003–2011, *J. Exposure Sci. Environ. Epidemiol.* 26 (4) (2016) 377.

- 639 [43] Y. Zhan, Y. Luo, X. Deng, K. Zhang, M. Zhang, M. L. Grieneisen, B. Di,
640 Satellite-based estimates of daily no₂ exposure in china using hybrid ran-
641 dom forest and spatiotemporal kriging model, *Environmental science &*
642 *technology* 52 (7) (2018) 4180–4189.
- 643 [44] A. C. Just, R. O. Wright, J. Schwartz, B. A. Coull, A. A. Baccarelli,
644 M. M. Tellez-Rojo, E. Moody, Y. Wang, A. Lyapustin, I. Kloog, Using
645 high-resolution satellite aerosol optical depth to estimate daily pm_{2.5} geo-
646 graphical distribution in mexico city, *Environmental science & technology*
647 49 (14) (2015) 8576–8584.
- 648 [45] U.S. Department of Commerce National Oceanic and Atmospheric Admin-
649 istration, Rapid update cycle (2018).
650 URL <https://ruc.noaa.gov/>
- 651 [46] C. Wiedinmyer, S. Akagi, R. J. Yokelson, L. Emmons, J. Al-Saadi, J. Or-
652 lando, A. Soja, The fire inventory from near (finn): a high resolution global
653 model to estimate the emissions from open burning, *Geosci. Model Dev.*
654 4 (3) (2011) 625.
- 655 [47] J. Fry, G. Z. Xian, S. Jin, J. Dewitz, C. G. Homer, L. Yang, C. A. Barnes,
656 N. D. Herold, J. D. Wickham, Completion of the 2006 national land cover
657 database for the conterminous united states, *Photogramm. Eng. Remote*
658 *Sens.* 77 (9) (2011) 858–864.
- 659 [48] P. F. Levelt, G. H. van den Oord, M. R. Dobber, A. Malkki, H. Visser,
660 J. de Vries, P. Stammes, J. O. Lundell, H. Saari, The ozone monitoring
661 instrument, *IEEE Trans. Geosci. Remote Sens.* 44 (5) (2006) 1093–1101.
- 662 [49] J. G. Powers, J. B. Klemp, W. C. Skamarock, C. A. Davis, J. Dudhia, D. O.
663 Gill, J. L. Coen, D. J. Gochis, R. Ahmadov, S. E. Peckham, G. A. Grell,
664 The weather research and forecasting model: Overview, system efforts, and
665 future directions, *Bull. Am. Meteorol. Soc.* 98 (8) (2017) 1717–1737.

- 666 [50] G. Pfister, D. Parrish, H. Worden, L. Emmons, D. Edwards, C. Wiedin-
667 myer, G. Diskin, G. Huey, S. Oltmans, V. Thouret, A. Weinheimer, Charac-
668 terizing summertime chemical boundary conditions for airmasses entering
669 the us west coast, *Atmos. Chem. Phys.* 11 (4) (2011) 1769–1790.
- 670 [51] M. Kuhn, C. Wing, S. Weston, A. Williams, C. Keefer, A. Engelhardt,
671 T. Cooper, Z. Mayer, B. Kenkel, the R Core Team, M. Benesty, R. Lescar-
672 beau, A. Ziem, L. Scrucca, Y. Tang, C. Candan, T. Hunt, caret: Classifi-
673 cation and Regression Training, r package version 6.0-80 (2018).
674 URL <https://CRAN.R-project.org/package=caret>
- 675 [52] J. Allaire, F. Chollet, keras: R Interface to 'Keras', r package version
676 2.2.0.9001.
677 URL <https://keras.rstudio.com>
- 678 [53] H. Singh, C. Cai, A. Kaduwela, A. Weinheimer, A. Wisthaler, Interactions
679 of fire emissions and urban pollution over california: Ozone formation and
680 air quality simulations, *Atmospheric environment* 56 (2012) 45–51.
- 681 [54] L. Breiman, Random forests, *Mach. Learn.* 45 (1) (2001) 5–32.