

1     Using Machine Learning to estimate the Impact of  
2     Ports and Cruise ship traffic on urban air quality: the  
3     case of Barcelona

4             Alexandre Fabregat<sup>a,\*</sup>, Lluís Vázquez<sup>a</sup>, Anton Vernet<sup>a</sup>

5             <sup>a</sup>*Department of Mechanical Engineering, Universitat Rovira i Virgili, Av. Països*  
6             *Catalans, 26, 43007, Tarragona (Spain)*

---

7     **Abstract**

Maritime activity is known to increase pollutant concentration levels in neighboring cities. In major touristic destinations, the singular need of cruise liners to keep supplying energy to on-board services and amenities while docked, has raised concerns about this industry contribution to pollutant emissions. To estimate the impact of port activities and that exclusively due to cruises, classical approaches would rely on atmospheric dispersion models. Although these tools retain the underlying physics, lack of details on background flow state and emission inventories limits their predictive capabilities. Using historical data on pollutant concentration, meteorology and traffic intensity at specific locations across the city of Barcelona, it was found that predictions of local pollutant concentration by the present Machine Learning tool are more accurate than those provided by the CALIOPE-Urban-v1.0 in our test cases. Estimated air quality impact due to cruise ships is shown to be limited in comparison to overall Port effects.

8     *Keywords:* Urban air pollution, Cruise ships, Generalized Boosted  
9     Regression Models, Machine Learning

10 **Highlights**

- 11     • Machine Learning is proposed as an alternative to classical dispersion  
12       models.
- 13     • Working dataset build from pollutant concentration, weather and traf-  
14       fic intensity.
- 15     • ML local predictions found to be more accurate than those from stan-  
16       dard approaches.
- 17     • Main features explaining pollutant concentration variability are iden-  
18       tified.
- 19     • Cruise ship activity impact on air quality of Barcelona metro area is  
20       quantified.

## 21 1. Introduction

22       Transportation represents a significant fraction of the total emissions of  
23 several chemical compounds that directly impact human health including  
24 Carbon Monoxide, Nitrogen and Sulfur Oxides, particulate matter, tropo-  
25 spheric Ozone and Volatile Organic Compounds (VOC) to mention the most  
26 relevant [1, 2]. Maritime traffic contribution to global emissions, both freight  
27 and passengers, has continuously grown over the last decades [3, 4, 5, 6]. De-  
28 spite notable efforts to reduce its impact by, for instance, limiting the sulphur  
29 content in marine fuel [7], globally, maritime transport is still characterized  
30 by consuming low quality fuel in comparison to road and air traffic [8]. At  
31 port, shipping activities have been found to contribute to increased levels of  
32 air pollutants in neighboring urban areas [9, 10].

33       Cruise shipping has raised special concerns because of its particular mode  
34 of operation characterized by the continuous supply of fossil fuel energy to  
35 the various on-board services and amenities during docking ('hotelling') [11].  
36 The average energy consumption used in manoeuvring cruise liners has been  
37 estimated to be only 5% of that used at dock to provide light, heating, venti-  
38 lation, air-conditioning, cold storage, cooking and other services estimated at  
39 2.5 MW per average cruise ship [12]. On one hand, previous efforts to quan-  
40 tify the impact of ports on air quality have been directed to characterize  
41 the pollutant emissions of docked vessels (number and location of sources,  
42 duration and rate, composition...) and use this information to feed more  
43 or less complex atmospheric dispersion models to predict the air quality in  
44 neighboring areas [12, 13, 9, 14]. Other studies analyzed pollutant concen-  
45 tration measurements to statistically infer connections between overall port

46 activity and air quality levels [15, 16, 17]. Isolated impact of cruise ships and  
 47 its peculiar regime of sustained emissions during hotelling on major touris-  
 48 tic destinations and areas of ecological value has also been the subject of  
 49 research [15, 18, 19, 14, 13].

In general, the transport of a chemical species  $i$  is governed by the convection-diffusion-reaction equation that can be written as:

$$\frac{\partial \rho Y_i}{\partial t} + \nabla \cdot (\rho \vec{u} Y_i) = -\nabla \cdot \vec{J}_i + R_i + S_i \quad (1)$$

50 where  $\rho$  is the background (mixture) fluid density,  $t$  is time,  $\vec{u}$  is the mixture  
 51 velocity field,  $Y_i$  is the local mass fraction of chemical species  $i$ ,  $\vec{J}_i$  is the  
 52 mass diffusion flux,  $R_i$  accounts for any production/consumption of  $i$  due to  
 53 chemical reactions and  $S_i$  represents any other source/sink of  $i$ .

54 Full solutions of Eq. (1) are usually unavailable for, at least, two reasons.  
 55 On one hand, local atmospheric velocity field  $\vec{u}$  responsible for advecting the  
 56 different chemical species in the turbulent planetary boundary layer is typi-  
 57 cally not well specified. Secondly, rates of emission, consumption by chemical  
 58 reaction and surface deposition in pollutant inventories are, in many cases,  
 59 approximations of varying accuracy. Several orders of magnitude smaller  
 60 than the convective transport counterpart, diffusive contribution is usually  
 61 neglected.

62 Classical approaches have sought for  $Y_i(\vec{x}, t)$  by numerically integrating  
 63 some approximated form of Eq. (1) [20, 21, 22]. In Gaussian Plume Mod-  
 64 els, one of the most commonly used approximation [23], wind velocity and  
 65 direction are assumed constant and turbulent transport is parametrized us-  
 66 ing horizontal and vertical standard deviations of the emission distribution

67 [1]. Other approaches to urban pollutant dispersion modelling involve the  
68 spatial and temporal discretization of some convenient form of the chemi-  
69 cal species transport equation Eq. (1) over a computational domain covering  
70 the geographical region and time period of interest. Instantaneous veloc-  
71 ity and concentration are approximated by some form of averaged/filtered  
72 fields resulting from modelling the turbulent contribution to the transport of  
73 momentum, mass and heat. Common methodologies involve embedded com-  
74 putational meshes with a coarse far field level where boundary conditions are  
75 imposed and progressively finer grids where detail is needed. Although this  
76 approach retains part of the physics that governs atmospheric dispersion, it  
77 also has important inherent sources of error due to the impossibility of recon-  
78 structing the full turbulent hydrodynamic field and the lack of detailed char-  
79 acterization of each pollutant emission and chemical transformation rates.

80 The approach presented here uses Supervised Machine Learning (ML)  
81 techniques to elucidate the relation between the local concentration of a given  
82 pollutant (the target or response) and several relevant variables for pollutant  
83 transport including meteorology and traffic intensity (the features or predic-  
84 tors). Use of Machine Learning algorithms to exploit the vast amounts of  
85 data retrieved from city-wide monitor networks has grown in recent years  
86 [24, 25, 26, 27]. While classical approaches retain the underlying physics of  
87 pollutant dispersion and are well suited for assessing the impact of different  
88 potential scenarios on local air quality (by, for example, considering differ-  
89 ent spatial locations and emission rates for individual pollutant sources),  
90 when robust and accurate estimations of the functionality between the fea-  
91 tures and the response exist, ML methodologies can accurately predict the

92 isolated effect of each predictor on pollutant concentration levels regardless  
93 of the availability of precise descriptions of the background flow state and  
94 pollutant inventories.

95 The model presented in this work is especially intended to estimate the  
96 overall impact of port activity and that exclusively due to cruise ships on the  
97 air quality of a major coastal touristic destination. The city of Barcelona has  
98 been chosen as a benchmark because it is one of the major travel destination  
99 with a total of  $\sim 9$  million tourists in 2018 [28, 29] and its Port impact on the  
100 metropolitan area has already been the subject of significant research efforts  
101 [30, 31, 32].

102 This work is organized as follows: Sec. 2 describes the methods and ma-  
103 terials involved in obtaining the working dataset used to train and tune the  
104 air quality predictive model. Sec. 3 presents the main results and quantifies  
105 the contribution of the Port of Barcelona activity to increased pollutant con-  
106 centration for several chemical species and locations across the metropolitan  
107 area. Sec. 4 discusses the hypothesis, assumptions and potential future im-  
108 provements to the current methodology and Sec. 5 summarizes the major  
109 findings.

## 110 **2. Methods**

### 111 *2.1. Overview*

In order to analyze the impact of the port activity and cruise ship traffic,  
we present and discuss an implementation of Supervised Machine Learning  
techniques for prediction of air quality levels in urban environments. The

specific goal is to obtain, for each pollutant species and selected spatial location, accurate estimations of the function  $f$  that captures the dependence between the local pollutant concentration — or *response*  $Y$  — and a set of  $n$  variables — or *predictors*  $X_j$  — this response depends upon, i.e.,

$$Y = f(X_j) + \varepsilon = \hat{Y} + \varepsilon, \quad j = 1 \dots n, \quad (2)$$

112 where  $\hat{Y}$  is the predicted concentration and  $\varepsilon$  is the error associated to un-  
113 explained variability in  $Y$ .

114 By inspection of Eq. (1), predictors must include information on both  
115 local atmospheric state and source emission rates of each pollutant. Precise  
116 estimations of pollutant concentration will be available as long as data on  $X_j$   
117 and  $Y$  allow to elucidate the function  $f$  that explains most of the variability  
118 in  $Y$  and minimizes  $\varepsilon$ , this is, the difference between observed and predicted  
119 values.

## 120 *2.2. Dataset description*

121 The working dataset used in this work and described in this section  
122 contains pollutant concentration measurements, weather measurements and  
123 road, air and maritime traffic data.

### 124 *2.2.1. Pollutant concentration*

125 The response dataset is comprised of hourly concentrations measurements  
126 of Nitric Oxide (NO), Nitrogen Dioxide (NO<sub>2</sub>), Total Nitrogen Oxides (NO<sub>x</sub>),  
127 Sulfur Dioxide (SO<sub>2</sub>), Ozone (O<sub>3</sub>), Particulate Matter 10 μm or less in diam-  
128 eter (PM<sub>10</sub>) and Carbon Monoxide (CO). The data is retrieved from  $M = 8$   
129 pollutant stations belonging to the Air Pollution Monitoring and Forecast-  
130 ing Network (XVPCA) of the Catalan Government [33] that are distributed

131 across the metropolitan area of Barcelona as shown in Fig. 1 (dark red mark-  
132 ers). Height and type of station is shown in Tab. 1. All pollutant concen-  
133 tration is in  $\mu\text{g m}^{-3}$  except CO that is in  $\text{mg m}^{-3}$ . Note that the relation  
134 between concentration and mass fraction of species  $i$  can be expressed as  
135  $c_i = Y_i \rho / M_i$  where  $M_i$  is the molar weight of  $i$ .

### 136 *2.2.2. Weather*

137 The weather dataset includes measurements on wind velocity ( $\text{km h}^{-1}$ )  
138 and direction ( $^\circ$ ), relative humidity (%), atmospheric pressure (hPa), pre-  
139 cipitation (mm), temperature ( $^\circ\text{C}$ ) and solar irradiance ( $\text{W m}^{-2}$ ) with a 30-  
140 minutes sampling rate from  $N = 11$  different meteorological stations across  
141 the metropolitan area of Barcelona and surroundings (the location of the  
142 closest stations to the metropolitan area are shown in Fig. 1 in blue). These  
143 data have been obtained from the network of automatic weather stations  
144 (XEMA) of the Meteorological Service of Catalonia [34].

145 This predictor set is extended with daily measurements of cloud cover  
146 from a single station located at the Observatori Fabra station (see Fig. 1)  
147 obtained from the European Climate Assessment and Dataset (ECA&D) [35].

### 148 *2.2.3. Traffic*

149 Pollutant emissions from combustion of maritime fossil fuels are incor-  
150 porated into the model by parsing traffic log files of the Port of Barcelona  
151 [36] (magenta marker in Fig. 1). The available data includes arrival and  
152 departure time stamps for each vessel and ship type (used to identify cruise  
153 liners) and size (length and beam). This data has been used to generate two  
154 predictors for the hourly number and median size of total vessels and cruise

155 liners.

156 Predictors for air traffic activity consist of hourly number of arrival and  
157 departures from the Airport of Barcelona [37] (green marker in Fig. 1).

158 Hourly predictors for road traffic intensity expressed as an index ranging  
159 from 1 (fluid) to 6 (total congestion) have been obtained from the extensive  
160 network of traffic measurement points across the city [38] (orange markers in  
161 Fig. 1).

162 Distributions of the number of total vessels and cruise ships simultane-  
163 ously docked at the Port of Barcelona over the time span considered in this  
164 work ranging from October 2017 to March 2020 (both included) as well as  
165 the distribution of flight operations at the Airport of Barcelona are shown in  
166 Appendix B.

### 167 *2.3. Preprocessing*

168 Weather and pollutant concentration measurements were joined by their  
169 respective observation time stamps into a synchronized single dataset. Daily  
170 measurements of cloud cover from a single weather station (Barcelona—  
171 Observatori Fabra) were assumed constant over the corresponding 24 hour  
172 period and over the entire metropolitan area. The working dataset was com-  
173 pleted by merging the three traffic data collections by their time stamps.

174 Estimations for road traffic intensity  $\sigma$  and every weather predictor  $\theta$  on  
175 each pollutant station  $i = \{1 \dots M\}$  were obtained by weighted interpolation  
176 of each observation. The weighting factor was set to the 4-th power of the  
177 distance  $d_{ij}$  between the pollutant station  $i$  and each of the  $j = \{1 \dots N\}$   
178 meteorological stations to ensure that closest observations contributes the  
179 most to the interpolated weather predictor. Analogously, the distance  $d_{ik}$

180 between the pollutant station  $i$  and each of the  $k = \{1 \dots Q\}$  road transit  
 181 measurement points were used in the spatial interpolation of traffic intensity.  
 182 Specifically,

$$\theta_i = \sum_{j=1}^N \hat{\theta}_j \eta_{ij}, \quad \sigma_i = \sum_{k=1}^Q \hat{\sigma}_k \eta_{ik}, \quad (3)$$

where  $\eta_{ij}$  is the distance weight from each pollutant station  $i$  to each meteorological station  $j$  and  $\eta_{ik}$  is the distance weight from each pollutant station  $i$  to each traffic measurement location  $k$ :

$$\eta_{ij} = \frac{s_{ij}}{\sum_{j=1}^N s_{ij}}, \quad \eta_{ik} = \frac{s_{ik}}{\sum_{k=1}^Q s_{ik}}, \quad (4)$$

where

$$s_{ij} = \frac{1}{d_{ij}^4}, \quad s_{ik} = \frac{1}{d_{ik}^4} \quad (5)$$

183 The wind direction predictor is expressed in terms of its alignment with  
 184 respect to the direction connecting each pollutant station  $i$  and the two major  
 185 transportation infrastructures likely to affect the air quality in the metropoli-  
 186 tan area, the Port and the Airport of Barcelona.

Given each pollutant station bearing angle with respect to the Port  $\alpha_{iP}$  and the Airport  $\alpha_{iA}$  and  $\beta_i$  being the direction the wind blows to, alignment with respect to the Port and Airport can be computed as:

$$a_{iP} = \left| \frac{|\alpha_{iP} - \beta_i|}{180} - 1 \right|, \quad a_{iA} = \left| \frac{|\alpha_{iA} - \beta_i|}{180} - 1 \right|, \quad (6)$$

187 where  $\alpha_{iP}$ ,  $\alpha_{iA}$  and  $\beta_i$  are expressed in degrees with the origin pointing to the  
 188 North and increasing in the clockwise (east) direction. The values of  $a_{iP}$  and  
 189  $a_{iA}$  are constrained in the range  $0 \leq a_i \leq 1$  where 0 indicates wind blowing  
 190 against the  $i$ -th station as seen from the corresponding infrastructure and 1  
 191 indicates complete alignment.

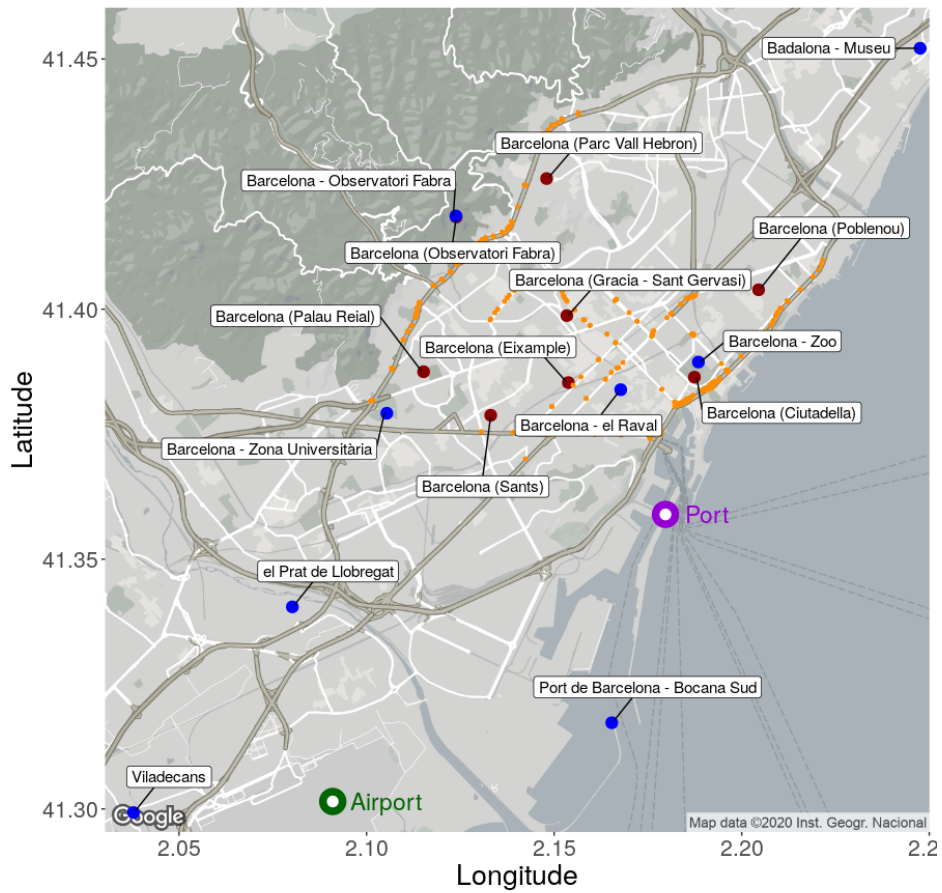


Figure 1: Map of Barcelona showing the location of pollutant (dark red) and meteorological (blue) measurement stations and traffic intensity measurements (orange). Cloud cover data comes from a single station (Observatori Fabra). Representative locations of the Port and Airport of Barcelona are shown in magenta and green respectively.

192 *2.4. Working dataset*

193 The working dataset is only comprised of complete observations, this is,  
 194 all predictors for a given hourly time stamps are known. The temporal span of  
 195 the working dataset ranges from October 2017 to March 2020 both included.  
 196 The list of the 25 predictors used in this work including a description and the  
 197 variable name is shown in Tab. 2. The data completeness for each pollutant  
 198 and station defined as the ratio between the actual number of observations  
 199 and the theoretical number of hourly observations over the 30 months period  
 200 is shown in Tab. 1. Note that while nitrogen oxides concentration data exists  
 201 for all stations, other pollutants are not available for all measurement points.

Station name	Data completeness (%)							Height (m)	Type
	CO	NO	NO <sub>2</sub>	NO <sub>x</sub>	O <sub>3</sub>	PM <sub>10</sub>	SO <sub>2</sub>		
Ciutadella	—	80	80	80	79	—	—	7	BU
Eixample	79	80	80	80	79	79	80	26	TU
Gràcia-Sant Gervasi	80	79	79	79	80	68	80	57	TU
Observatori Fabra	—	63	63	63	61	43	—	415	BSU
Palau Reial	72	74	74	74	71	67	74	81	BU
Parc Vall Hebron	75	78	78	78	78	75	79	136	BU
Poblenou	—	64	64	64	—	62	—	3	BU
Sants	—	79	79	79	—	—	—	35	BU

Table 1: Pollutant station data completeness (%) for each pollutant, height and type (‘BU’ = Background Urban, ‘BSU’ = Background Suburban, ‘TU’ = Traffic urban).

202 To illustrate the data contained in the working dataset, Fig. 2 shows the  
 203 temporal evolution of some variables (including the response for NO<sub>x</sub>) for  
 204 the Barcelona—Eixample station over March 2019.

Predictor description	Variable
Hourly wind alignment between Port and pollutant station	port_align
Hourly wind velocity	wind_velocity
Daily cloud cover	cloud_cover
Hourly wind alignment between Airport and pollutant station	airport_align
Hourly precipitation	precipitation
Hourly temperature	temperature
Hourly atmospheric pressure	pressure
Hourly solar irradiance	irradiance
Hourly relative humidity	rel_humidity
Hourly traffic intensity	road_traffic
Hourly total number of vessels	number_vessels
Hourly number of cruise ships	number_cruises
Median length of hourly total vessels	vessel_length
Median length of hourly cruise ships	cruise_length
Day is Monday	monday
Day is Tuesday	tuesday
Day is Wednesday	wednesday
Day is Thursday	thursday
Day is Friday	friday
Day is Saturday	saturday
Day is Sunday	sunday
Day number of the year	day_number
Year number	year
Hour of the day	hour
Hourly number of flights into/from the Airport of Barcelona	air_traffic

Table 2: List of predictors (or features) and corresponding variable names.

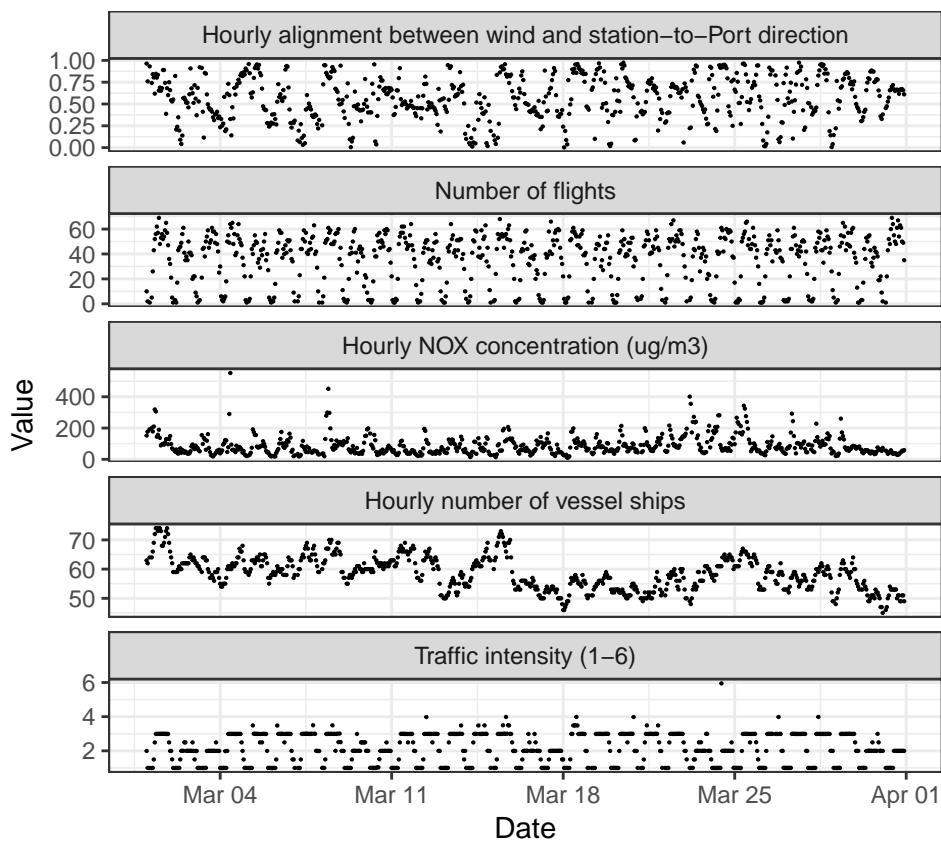


Figure 2: Temporal evolution of hourly wind alignment with respect to the Barcelona—Eixample station, number of flights, hourly concentration of  $\text{NO}_x$ , Hourly number of total vessels and traffic intensity over March 2019.

205 *2.5. ML procedure*

206 The R package *mlr 2.17.1* [39] has been used for data pre-processing,  
 207 resampling, hyperparameter tuning and post-processing. Regarding the pre-  
 208 processing, all continuous features have been normalized by subtracting the  
 209 average value and dividing by the standard deviation for each pollutant  
 210 species and measurement station. Categorical variables have been converted

211 to continuous features using a *1-of-1* method.

212 In this work, six different learning algorithms (or learners) widely used  
213 in ML applications have been tested, namely Gradient Boosting Machine  
214 or GBM (*gbm*), Support Vector Machines or SVM (*ksvm*), Random Forest  
215 (*h2o.radnomForest*), Multivariate Adaptive Regression Splines (*earth*), Feed-  
216 Forward Multilayer Artificial Neural Network (*h2o.deeplearning*) and Feed-  
217 Forward Neural Networks for Multinomial Log-linear models (*nnet*) where  
218 the name of the integrated regression model in the R package *mlr 2.17.1* [39]  
219 appears in parenthesis. A 5-fold cross validations resampling strategy has  
220 been used to assess the performance of each learner. Thus, the working data  
221 set described above has been repeatedly split into a training and a test set.  
222 After optimizing each learner hyperparameters, the training set has been  
223 used to derive the corresponding model which performance has then been  
224 estimated using the testing set. The benchmark results in terms of mean  
225 squared error for each pollutant across all stations for each learner are shown  
226 in Appendix C. Results clearly indicate that, for this specific problem and  
227 working dataset, the Gradient Boosting Machine (GBM) outperforms the rest  
228 of the algorithms considered in this work. Optimal GBM hyperparameters  
229 were found to be around 5000 trees and an interaction depth of 25 for most  
230 of pollutants and measurement stations. All results in Sec. 3 have been  
231 obtained with the GBM methodology.

232 All computations have been performed on a Dell AIO 7770 workstation  
233 with a Intel(R) Core(TM) i9-9900 CPU @ 3.10GHz. Learner training and  
234 testing, hyperparameter tuning and resampling typical CPU time is 45 min-  
235 utes for each station and pollutant. Single point prediction takes less than 1

236 ms.

## 237 **3. Results**

### 238 *3.1. Validation*

239 To assess the predictive performance of the ML-based tool presented in  
240 Sec. 2, predicted concentration levels are compared with those reported by  
241 Benavides *et al.* [40] who used the CALIOPE-Urban v1.0 *physics-based*  
242 model to predict NO<sub>2</sub> concentration at three of the pollutant measurement  
243 stations listed in Tab. 1. CALIOPE-Urban v1.0 is the result of coupling  
244 CALIOPE [41] — an operational mesoscale air quality forecast system based  
245 on the HERMES (emissions) [42], WRF (meteorology) [43] and CMAQ  
246 (chemistry) [44] models — with the urban roadway dispersion model R-LINE  
247 [45].

248 The results shown in Tab. 3 indicate that the *ML-based* model outper-  
249 forms the CALIOPE-Urban v1.0 platform in every metric for all available  
250 stations. Definitions of Geometric Mean Bias (GeoMean), correlation coeffi-  
251 cient ( $R$ ), Mean Bias (MB) Root Mean Square Error (RMSE) and Fraction  
252 of predictions within a factor of 2 in observations (FAC2) can be found in  
253 the Appendix A.

### 254 *3.2. Feature importance*

255 The relative importance of each predictor has been determined using the  
256 entropy-based *information gain* [39] that returns an *importance weight* for  
257 each feature that maximizes the amount of captured variability in the re-  
258 sponse. The four most important features for each station are shown in

Metric	Eixample		Palau Reial		Grcia-Sant Gervasi	
	Benavides <i>et al.</i> [40]	Current	Benavides <i>et al.</i> [40]	Current	Benavides <i>et al.</i> [40]	Current
GeoMean	0.83	0.97	1.10	0.92	1.07	0.96
$R$	0.55	0.81	0.57	0.78	0.52	0.82
MB	8.57	-0.04	1.23	0.07	6.00	-0.10
RMSE	26.70	14.16	21.57	13.12	25.11	14.74
FAC2	0.86	0.97	0.73	0.85	0.79	0.94

Table 3: Comparison of statistics between Benavides *et. al* and the current study.

259 Fig. 3 in descending order from top to bottom rows. Results indicate that  
260 the day of the year is the the most important feature for almost all back-  
261 ground urban (in black) and the only background suburban (in red) stations.  
262 Although this predictor also explains most of the variability in traffic urban  
263 stations (in blue), this leading role is shared with the year number in the  
264 case of CO, the traffic intensity for NO<sub>x</sub> and the temperature for PM<sub>10</sub>. The  
265 hour of the day is the second most important feature for NO<sub>x</sub> and O<sub>3</sub> at most  
266 stations, with some exceptions for which this position is occupied by traffic  
267 intensity, wind speed and temperature. In the case of CO, the second most  
268 relevant predictor at the traffic urban stations is the traffic intensity while  
269 the year number dominates in background urban sites. This latter feature is  
270 also the second most important variable in the case of SO<sub>2</sub>. This suggests  
271 that, despite the relatively small number of years in the working dataset,  
272 the concentrations for these two species exhibit significant changes over the  
273 2017-2020 period. In contrast to the other species, the PM<sub>10</sub> concentration  
274 at most locations has been found to be controlled in second place by the  
275 temperature.

276 While the top two most important features were dominated by temporally-  
277 related predictors, namely the hour of the day, the day of the year and the  
278 year number, the third and fourth positions, with a much smaller impact  
279 on the concentration variability as discussed later, exhibit notable variability  
280 across stations and pollutants. Thus, carbon Monoxide concentration even-  
281 tually exhibits some dependence on the docked ship size. Similar behavior is  
282 exhibited by the sulphur dioxide.  $\text{NO}_x$ , mostly controlled by the day of the  
283 year and the hour of the day, has also been found to change with the intensity  
284 of the traffic, the wind velocity and the vessel size. Third and fourth features  
285 for Ozone include the wind velocity and alignment with respect to the Port  
286 and the temperature. Finally, the particulate matter, strongly dependent on  
287 the day of the year and the temperature, exhibits also some dependence on  
288 the traffic intensity, the hour of the day, the wind alignment with respect to  
289 the Port and the number of flights at the Airport.

290 The *importance gain* method for each station and pollutant are shown  
291 in Fig. 4. This metric has been normalized such that its sum across all  
292 predictors is 100. In general, the results indicate a consistent distribution of  
293 importance between stations and pollutants with most of the variability in  
294 concentration explained by the hour of the day, the day number of the year,  
295 the road traffic intensity and the temperature. On the other side, weekday  
296 predictors are found to be the least important features.

297 Notably though, the model reveals a few notable peculiarities. First, the  
298 results in the suburban station (Observatori Fabra) located at an altitude  
299 more than twice that of the second one in height, suggest that ozone con-  
300 centration is mostly controlled by the day of the year. These results are in

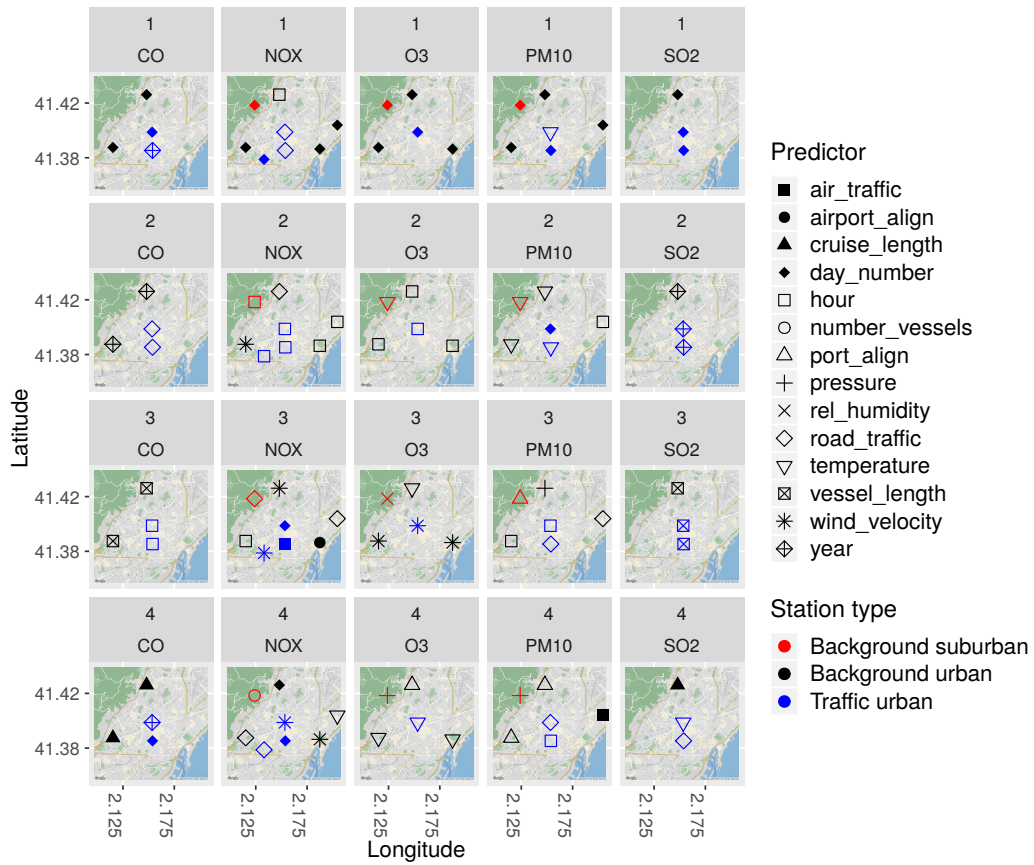


Figure 3: Map showing the top four most important predictors as measured by the entropy-based information gain for CO, NO<sub>x</sub>, O<sub>3</sub>, PM<sub>10</sub> and SO<sub>2</sub>. Colors indicate the type of pollutant station as shown in Fig. 1. Each symbol identifies the corresponding predictor as listed in Tab. 2.

301 agreement to the well-known seasonal behavior of tropospheric ozone with  
302 higher values in the summer months when atmospheric conditions favour O<sub>3</sub>  
303 formation.

304 As expected, traffic intensity scores high in importance in urban traffic  
305 stations (Eixample and Gràcia - Sant Gervasi) for nitrogen oxides and carbon  
306 monoxide. Contrarily, limited effects on these pollutant concentrations for  
307 this feature are found in background and suburban locations. The impact  
308 of traffic intensity on nitrogen oxides concentration at the Observatori Fabra  
309 station can be explained by the proximity of a major road traffic route (Ronda  
310 de Dalt). Wind velocity seems to affect also the nitrogen oxides level in  
311 background urban and suburban stations.

312 Wind velocity impact on concentration levels significantly varies across  
313 stations with Palau Reial and Sants exhibiting the largest effect on nitrogen  
314 oxide levels. Wind alignment with respect to the Airport and the Port are  
315 found to have similar importance with the latter exhibiting larger impact  
316 on PM<sub>10</sub> concentration at the two only ‘Traffic Urban’ stations, namely,  
317 Eixample and Gràcia - Sant Gervasi.

318 As happens for the two predictors for the number of total vessels and  
319 cruise liners, the ship size of the latter type is found to be the least relevant.  
320 The number of airport movements, the pressure and irradiance predictors  
321 exhibit modest importance and large cross-station and cross-pollutant vari-  
322 ability.

323 The overall least important features are those related to the day of the  
324 week. As shown, except Sunday — a holiday—, the workweek days do not  
325 show a great influence on pollutant concentration. Negligible importance has

326 also been found for hourly precipitation, cloud cover, relative humidity and  
 327 notably, the hourly number of docked cruise ships.

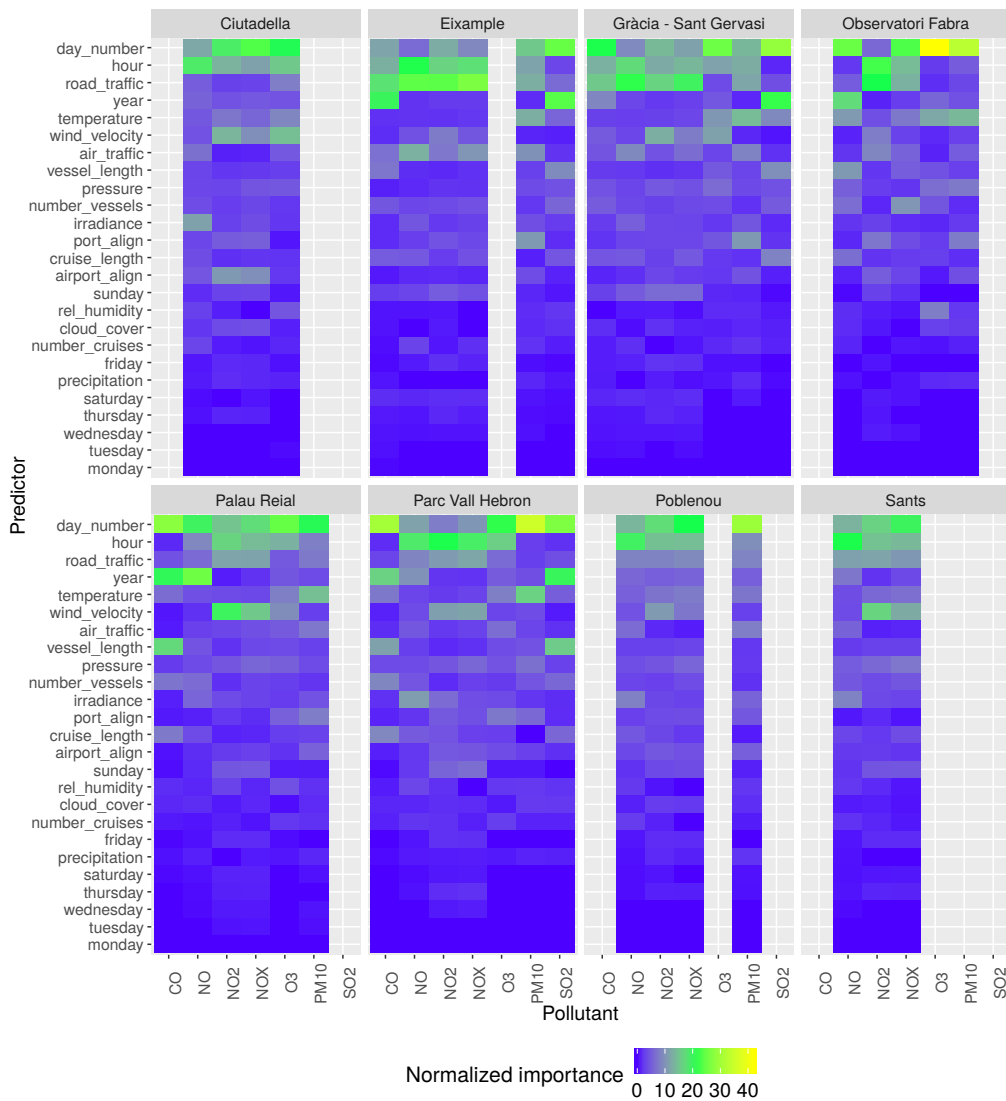


Figure 4: Normalized information gain attributable to each predictor in Tab. 2 as measured by the entropy-based information gain for each pollutant. Each panel corresponds to a different station (see Fig. 1).

328 *3.3. Overall Port impact*

329 Using the hourly number of total vessels as a proxy for the Port of  
330 Barcelona activity, the isolated impact of this predictor on the air quality  
331 at each pollutant station has been estimated by comparing the actual and  
332 the predicted concentration levels for several values of this predictor. Thus,  
333 the  $x$ -axis of Fig. 5 shows the difference between the total number of vessels  
334 used in the prediction and the actual value. The  $y$ -axis shows the average  
335 difference over all observations in the working dataset between the observed  
336 concentration and average prediction of the GBM learner (using a 3-fold re-  
337 sampling). For reference, mean and median concentration values for each  
338 station and pollutant in the working dataset are shown in Tab. 4.

339 As shown in Fig. 5, the excellent accuracy of the ML model results in  
340 negligible averaged difference between the actual and the predicted concen-  
341 tration ( $y \approx 0$ ) when the number of vessels used in the prediction coincides  
342 with the actual one, i.e.,  $x = 0$ .

343 In general, the model predicts that an increase in the activity in the Port  
344 of Barcelona as measured by the total vessel traffic leads to an increase in  
345 the local concentrations of carbon monoxide, nitrogen oxides and particulate  
346 matter. The overall trend indicates that the impact of the maritime traffic  
347 weakens as the distance from the Port increases. Thus, the impact of the  
348 number of vessels on the concentration levels in Ciutadella and Eixample  
349 stations is notably larger than that in Observatori Fabra and Palau Reial.

350 The largest impacts predicted by the model are found at the closest to  
351 the Port station (Eixample) where NO and NO<sub>2</sub> concentrations are expected  
352 to rise around 12% and 10% respectively in comparison to the mean value

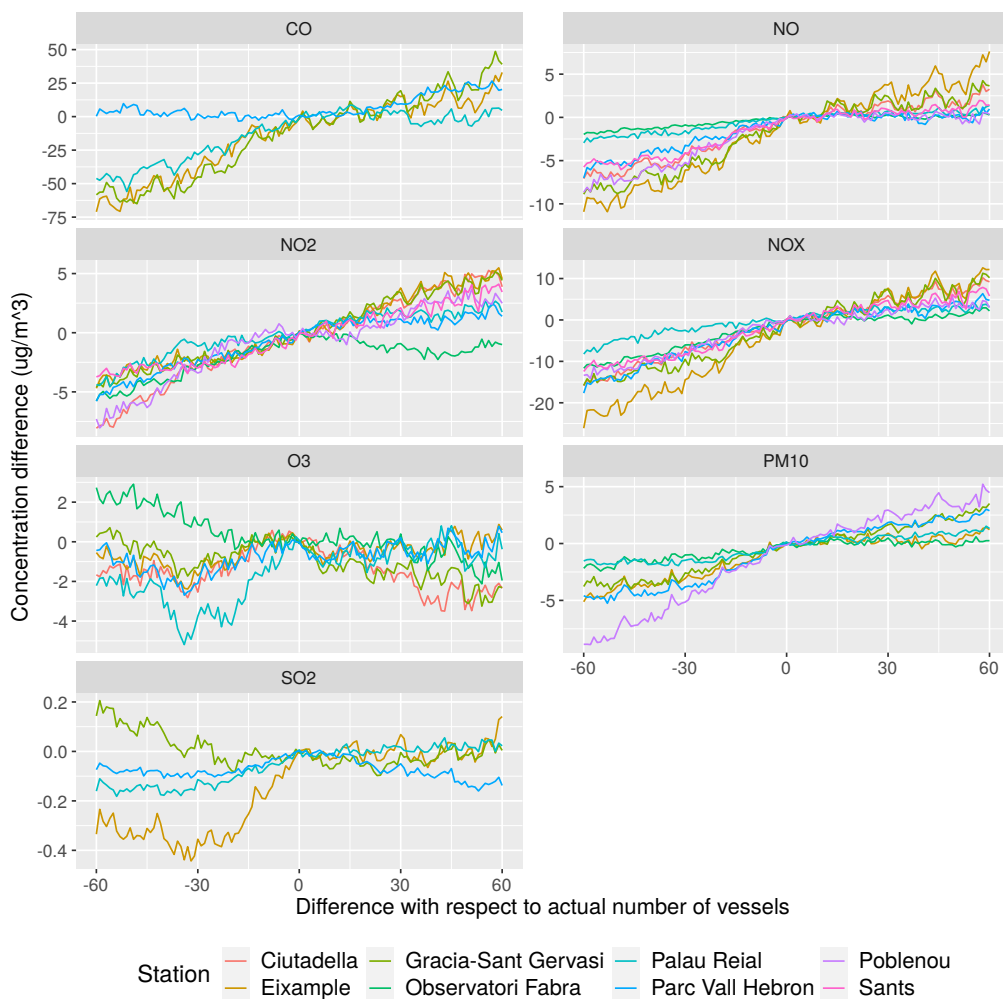


Figure 5: Predicted average concentration increase in  $\mu\text{g m}^{-3}$  due to overall Port traffic for each pollutant measurement station (see Fig. 1). Each panel corresponds to a pollutant species.

353 in Tab. 4. In contrast, the largest increases in CO (9%) and PM<sub>10</sub> (11%)  
 354 concentrations are found at slightly further stations (Grcia - Sant Gervasi and  
 355 Poblenu respectively). The signature of the Port activity on the ozone and  
 356 sulphur dioxide concentration is clearly weaker suggesting lack of correlation

357 between the number of vessels and their concentration levels.

Station	Metric	CO	NO	NO <sub>2</sub>	NO <sub>x</sub>	O <sub>3</sub>	PM <sub>10</sub>	SO <sub>2</sub>
Ciutadella	Mean	-	13.08	34.29	53.96	41.64	-	-
	Median	-	3	31	38	41	-	-
Eixample	Mean	410.0	32.47	52.16	101.65	36.51	26.61	1.79
	Median	300	19	49	79	35	24	1
Gràcia - Sant Gervasi	Mean	414.0	19.41	45.54	74.81	42.54	24.34	1.90
	Median	300	9	41	56	43	22	1
Observatori Fabra	Mean	-	2.75	11.71	16.17	82.85	16.57	-
	Median	-	1	9	11	82	14	-
Palau Reial	Mean	284.7	7.47	28.49	38.91	50.78	18.96	1.56
	Median	200	2	22	25	51	18	1
Parc Vall Hebron	Mean	294.0	8.17	29.15	41.36	54.14	21.42	1.42
	Median	300	3	23	27	53	19	1
Poblenou	Mean	-	15.07	38.01	60.84	-	28.11	-
	Median	-	5	34	43	-	24	-
Sants	Mean	-	9.03	32.65	46.24	-	-	-
	Median	-	3	27	32	-	-	-

Table 4: Mean and median concentration in  $\mu\text{g m}^{-3}$  for every pollutant and station.

### 358 3.4. Cruise ship impact

359 Following an analogous approach, the hourly number of liners has been  
 360 used as a proxy for cruise shipping activity in the Port of Barcelona. In  
 361 Fig. 6, the  $x$ -axis corresponds to the average difference between the actual  
 362 number of cruise ships and that used in the prediction while the  $y$ -axis shows  
 363 the corresponding difference between the observed and average predicted  
 364 concentration over all observations in the working dataset.

365 Of course, since cruise ships represent a relatively small fraction of the  
 366 total vessels, its isolated effect on the air quality of Barcelona is expected

367 to be less pronounced. Predictions suggest that the largest impact of cruise  
368 ships on nitrogen oxide concentration levels occurs at the two stations closest  
369 to the Port (Ciutadella and Eixample) with  $1.3 \mu\text{g m}^{-3}$  of  $\text{NO}_x$  for each  
370 additional cruise liner in Port. Impact on nitrogen oxides levels is predicted  
371 to decay as distance between the station and the Port increases.

372 In contrast to this marked variability across stations for  $\text{NO}_x$ , predicted  
373 concentrations of  $\text{PM}_{10}$  are found to vary in a consistent way across stations  
374 with an average slope of  $0.04 \mu\text{g m}^{-3}$  for each additional cruise liner.

375 CO is predicted to be insensitive when changes in the number of docked  
376 cruise liners is modest ( $\pm 3$  ships). For variations above  $\pm 5$  cruise ships,  
377 CO concentration levels start exhibiting significant variations at the closest  
378 stations to the Port. Farther locations are predicted to remain constant  
379 independently of the number of liners.

380 Cruise ship activity seems to have a negligible impact on sulfur dioxide  
381 with maximum variations in concentration less than 2.5% change with re-  
382 spect to the across-station mean value of  $1.67 \mu\text{g m}^{-3}$  (obtained from the  
383 four station mean values in Tab. 4). Predicted impact on ozone concen-  
384 trations is found to be negligible in most stations and negatively correlated  
385 in others, specially Observatori Fabra. The relatively large horizontal and  
386 vertical distance between the Port and this station suggest that, on overall,  
387 cruise ships do not affect the ozone levels across the metropolitan area.

388 The overall modest increase in pollutant concentration of the hourly num-  
389 ber of cruise liners and this predictor relatively small *importance gain* suggest  
390 that, in comparison to the overall effects of the Port, the impact of the cruise  
391 activity on the air quality of the city is very limited.

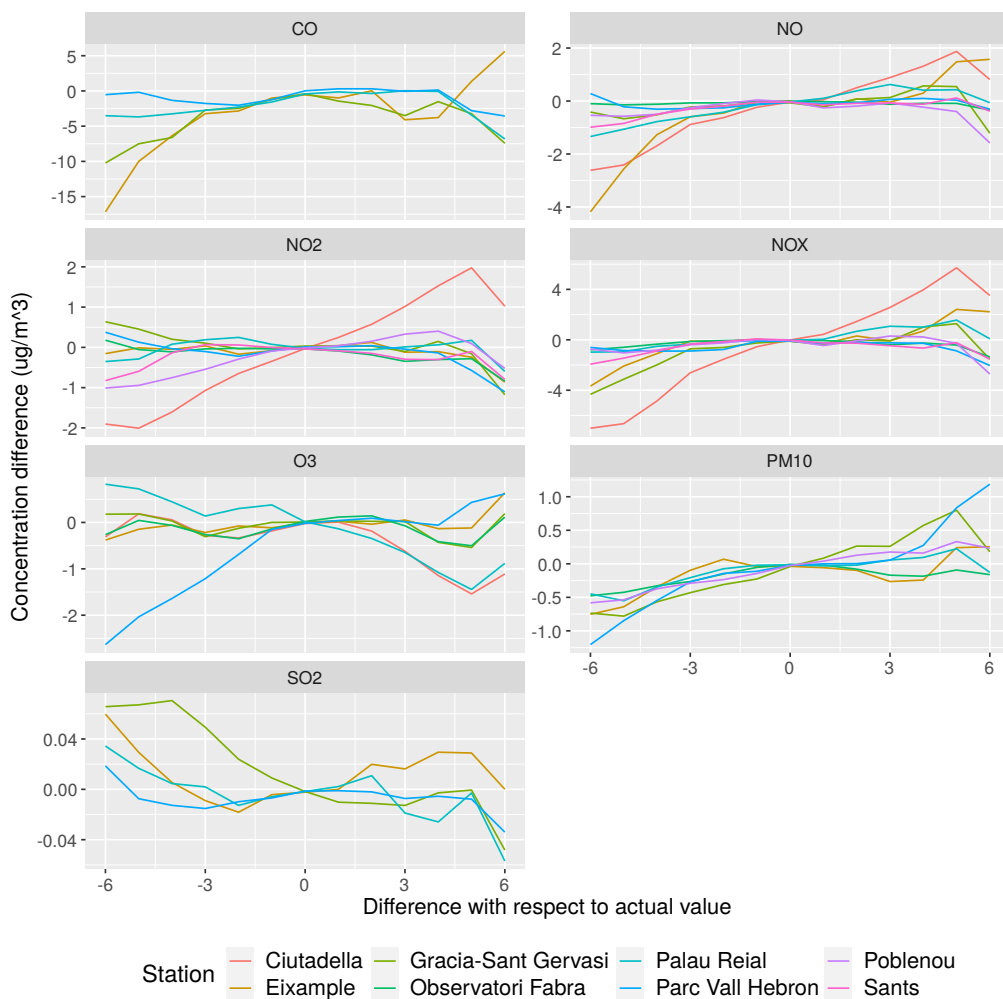


Figure 6: Predicted average concentration increase in  $\mu\text{g m}^{-3}$  due to cruise ship traffic for each pollutant measurement station (see Fig. 1). Each panel corresponds to a pollutant species.

392 **4. Discussion**

393 The availability and ease of access to large amounts of data from uncount-  
 394 able sources of scientific and technological relevance has made it possible to

395 use Machine Learning techniques to improve our understanding of processes  
396 and phenomena until now only accessible through classical approaches.

397 This paradigm is illustrated in this work where, provided enough data on  
398 both the atmospheric state and several proxies on pollutant inventories, it is  
399 shown that ML-based predictors are capable of outperforming traditional pol-  
400 lutant dispersion models. While traditional, physics-based approaches suffer  
401 from uncertainties on local weather and pollutant inventories, ML approaches  
402 are generally criticized by their *black box* nature offering small insight on the  
403 inner workings of the physical processes under study.

404 Besides these general considerations, and despite the excellent predictive  
405 capabilities of the tool developed here, this model has two main limitations  
406 in comparison to traditional approaches. On one hand, the relatively small  
407 number of stations reporting concentration levels across the metropolitan  
408 area of Barcelona hinders the potential of the model to be used as an op-  
409 erational tool capable of predicting air quality at any location across the  
410 region of interest. In other words, the current model lacks of spatial features.  
411 Therefore, air quality predictions are restricted to the locations and heights  
412 of each pollutant station for which an independent model has been derived.

413 Pollutant measurement campaigns with fixed or mobile equipment might  
414 provide in the future high spatial resolution information that make possible  
415 to use measurement location and altitude as predictors. This key feature  
416 would enable the transition of the current ML model into a fully-operational  
417 tool capable of providing accurate, fast and cheap predictions of air quality.

418 Second, in this work pollutant emissions are accounted via some proxy  
419 variables including the number of flights, vessels and road traffic congestion.

420 Although these are prime contributions to urban pollution, industrial and  
421 agricultural emission may also play a significant role in worsening the air  
422 quality of neighboring cities. Overall model accuracy could be potentially  
423 improved provided better characterizations of pollutant emission inventories  
424 in addition to transportation. Of course, the performance of the ML tool  
425 presented here is expected to increase as the size and quality of the time  
426 series datasets used to train the learners increase.

## 427 **5. Conclusions**

428 Weather and pollutant concentration measurements collected by monitor-  
429 ing networks across the metropolitan area of Barcelona have been combined  
430 with road, air and maritime traffic intensity data to train a Machine Learn-  
431 ing model to predict the air quality in the city. In terms of predictive capa-  
432 bilities of local  $\text{NO}_2$  concentration levels at several station locations across  
433 Barcelona, the resulting tool has been found to have better performance than  
434 the CALIOPE Urban v1.0 platform.

435 Using Mean Squared and Absolute Errors as metrics for predictive ac-  
436 curacy, the benchmark of different classical learners suggested that, for this  
437 specific dataset and across all stations and pollutants, the Gradient Boosting  
438 Machine (GBM) has been found to be the best performer.

439 Importance analysis based on information gain indicates that the three  
440 top features explaining pollutant level variability are the time of the year,  
441 the daytime and the road traffic intensity. The least important predictors  
442 include the workweek day, the precipitation, relative humidity and, notably,  
443 the number of cruise ships in Port.

444 Using the difference between observed and predicted pollutant concen-  
445 trations for several difference values with respect to the actual number of  
446 vessels, the ML tool has been used to estimate the contribution to worsened  
447 air quality of the overall Port activity and that exclusively due to cruise  
448 liners.

449 Results suggest that, close to the Port of Barcelona, the overall activity  
450 of this infrastructure leads to 12% and 10% increased mean concentration  
451 of NO and NO<sub>2</sub> respectively. Largest increases in CO and PM<sub>10</sub> have been  
452 estimated to be 9% and 11%. In comparison, the Port impact on SO<sub>2</sub> and  
453 O<sub>3</sub> concentration levels is predicted to be significantly weaker.

454 Using the hourly number of docked cruise ships as a proxy to estimate  
455 the isolated effect of this industry, predictions suggest that the impact of  
456 liners on the air quality is very limited in comparison to the overall Port  
457 contribution. Predictions suggest that the largest impact of cruise ships is  
458 1.3 µg m<sup>-3</sup> of NO<sub>x</sub> for each additional docked cruise ship in Port.

459 Besides the better accuracy in predicting pollutant concentration, the  
460 numerical simulations carried out with the ML-based model presented here  
461 required only a personal workstation with typical run-times in the order of  
462 hours. This computational costs are well below those typically required by  
463 state-of-the-art dispersion models.

## 464 **6. Acknowledgements**

465 This study was financially supported by the Spanish Ministry of Economy,  
466 Industry and Competitiveness — Research National Agency (under projects  
467 DPI2016-75791-C2-1-P and RTI2018-100907-A-I00), by FEDER funds and

468 by the Generalitat de Catalunya — AGAUR (under project 2017 SGR 01234).  
469 Authors would like to thank the Government of Catalonia for granting ac-  
470 cess to atmospheric pollutant data from the XVPCA network and the Servei  
471 Meteorològic de Catalunya for weather measurements.

472 [1] J. Seinfeld, S. Pandis, Atmospheric Chemistry and Physics: From Air  
473 Pollution to Climate Change, 1998.

474 [2] T. R. Walker, O. Adebambo, M. C. Del Aguila Feijoo, E. Elhaimer,  
475 T. Hossain, S. Johnston Edwards, C. E. Morrison, J. Romo, N. Sharma,  
476 S. Taylor, S. Zomorodi, Chapter 27 - Environmental Effects of Marine  
477 Transportation, in: C. Sheppard (Ed.), World Seas: an Environmental  
478 Evaluation (Second Edition), Academic Press, second edition edn., 505  
479 – 530, 2019.

480 [3] S. Sorte, V. Rodrigues, C. Borrego, A. Monteiro, Impact of harbour  
481 activities on local air quality: A review, Environmental Pollution 257  
482 (2019) 113542.

483 [4] M. Viana, P. Hammingh, A. Colette, X. Querol, B. Degraeuwe,  
484 de Vlieger I., J. van Aardenne, Impact of maritime transport emissions  
485 on coastal air quality in Europe, Atmospheric Environment 90 (2014)  
486 96 – 105.

487 [5] V. A. Profillidis, G. N. Botzoris, Chapter 2 - Evolution and Trends of  
488 Transport Demand, in: V. Profillidis, G. Botzoris (Eds.), Modeling of  
489 Transport Demand, Elsevier, 47 – 87, 2019.

- 490 [6] M. Zhao, Y. Zhang, W. Ma, Q. Fu, X. Yang, C. Li, B. Zhou, Q. Yu,  
491 L. Chen, Characteristics and ship traffic source identification of air pol-  
492 lutants in China's largest port, *Atmospheric Environment* 64 (2013)  
493 277286, doi:\bibinfo{doi}{10.1016/j.atmosenv.2012.10.007}.
- 494 [7] C. Schembari, F. Cavalli, E. Cuccia, J. Hjorth, G. Calzolari, N. Perez,  
495 J. Pey, P. Prati, F. Raes, Impact of a European directive on ship emis-  
496 sions on air quality in Mediterranean harbours, *Atmospheric Environ-*  
497 *ment* 61 (2012) 661669.
- 498 [8] D. Lack, J. Corbett, Black carbon from ships: a review of the effects of  
499 ship speed, fuel quality and exhaust gas scrubbing, *Atmospheric Chem-*  
500 *istry & Physics Discussions* 12 (2011) 3509–3554.
- 501 [9] L. Fileni, E. Mancinelli, M. Morichetti, G. Passerini, U. Rizza, S. Virgili,  
502 Air Pollution in Ancone Harbour, Italy, in: *Maritime Transport 2019*,  
503 199–208, 2019.
- 504 [10] J. Isakson, T. Persson, E. Lindgren, Identification and assessment of  
505 ship emissions and their effects in the harbour of Göteborg, Sweden,  
506 *Atmospheric Environment* 35 (2001) 3659–3666, doi:\bibinfo{doi}{10.  
507 1016/S1352-2310(00)00528-8}.
- 508 [11] E. Merico, A. Donateo, A. Gambaro, D. Cesari, E. Gregoris, E. Barbaro,  
509 A. Dinoi, G. Giovanelli, S. Masieri, D. Contini, Influence of in-port ships  
510 emissions to gaseous atmospheric pollutants and to particulate matter  
511 of different sizes in a Mediterranean harbour in Italy, *Atmospheric En-*  
512 *vironment* 139.

- 513 [12] H. Saxe, T. Larsen, Air pollution from ships in three Danish ports,  
514 Atmospheric Environment 38 (2004) 4057–4067.
- 515 [13] K. Poplawski, E. Setton, B. McEwen, D. Hrebenyk, M. Graham,  
516 C. Keller, Impact of cruise ship emissions in Victoria, BC, Canada,  
517 Atmospheric Environment - Atmos Environ 45 (2011) 824–833.
- 518 [14] F. Murena, M. Luigia, F. Quaranta, D. Toscano, Impact on air quality  
519 of cruise ship emissions in Naples, Italy, Atmospheric Environment 187.
- 520 [15] S. Eckhardt, O. Hermansen, H. Grythe, M. Fiebig, M. Cassiani,  
521 A. Baecklund, A. Stohl, The influence of cruise ship emissions on air  
522 pollution in Svalbard — A harbinger of a more polluted Arctic?, Atmo-  
523 spheric Chemistry & Physics Discussions 13 (2013) 3071–3093.
- 524 [16] A. Donateo, E. Gregoris, A. Gambaro, E. Merico, R. Giua, A. No-  
525 cioni, D. Contini, Contribution of harbour activities and ship traffic to  
526 PM<sub>2.5</sub>, particle number concentrations and PAHs in a port city of the  
527 Mediterranean Sea (Italy), Environmental science and pollution research  
528 international 21.
- 529 [17] D. Contini, A. Gambaro, F. Belosi, S. Pieri, W. Cairns, A. Donateo,  
530 E. Zanutto, M. Citron, The direct influence of ship traffic on atmo-  
531 spheric PM<sub>2.5</sub>, PM<sub>10</sub> and PAH in Venice, Journal of Environmental  
532 Management 92 (2011) 2119–2129.
- 533 [18] A. Maragkogianni, S. Papaefthimiou, Evaluating the social cost of cruise  
534 ships air emissions in major ports of Greece, Transportation Research  
535 Part D: Transport and Environment 36 (2015) 10 – 17.

- 536 [19] J. Celic, S. Valcic, M. Bistrović, Air pollution from cruise ships, Pro-  
537 ceedings Elmar - International Symposium Electronics in Marine (2014)  
538 75–78.
- 539 [20] L. T. Fan, Y. Horie, H. J. Paulus, Review of atmospheric dispersion and  
540 urban air pollution models, CRC Critical Reviews in Environmental  
541 Control 2 (1-4) (1972) 431–457.
- 542 [21] A. Leelossy, F. Molnár, F. Izsák, A. Havasi, I. Lagzi, R. Mészáros, Dis-  
543 persion modeling of air pollutants in the atmosphere: a review, Central  
544 European Journal of Geosciences 6 (2014) 257–278, doi:\bibinfo{doi}  
545 {10.2478/s13533-012-0188-6}.
- 546 [22] M. Lateb, R. N. Meroney, M. Yataghene, H. Fellouah, F. Saleh, M. C.  
547 Boufadel, On the use of numerical modelling for near-field pollutant  
548 dispersion in urban environments — A review, Environmental Pollution  
549 208 (2016) 271 – 283, ISSN 0269-7491, doi:\bibinfo{doi}{https://doi.  
550 org/10.1016/j.envpol.2015.07.039}, URL [http://www.sciencedirect.  
551 com/science/article/pii/S0269749115003723](http://www.sciencedirect.com/science/article/pii/S0269749115003723), special Issue: Urban  
552 Health and Wellbeing.
- 553 [23] C. Hood, I. MacKenzie, J. Stocker, K. Johnson, D. Carruthers,  
554 M. Vieno, R. Doherty, Air quality simulations for London using  
555 a coupled regional-to-local modelling system, Atmospheric Chem-  
556 istry and Physics 18 (2018) 11221–11245, doi:\bibinfo{doi}{10.5194/  
557 acp-18-11221-2018}.
- 558 [24] H. Relvas, A. Miranda, An urban air quality modeling system to support

- 559 decision-making: design and implementation, *Air Quality, Atmosphere*  
560 & *Health* 11, doi:\bibinfo{doi}{10.1007/s11869-018-0587-z}.
- 561 [25] J. Song, K. Han, Deep-MAPS: Machine Learning based Mobile Air Pol-  
562 lution Sensing, 2019.
- 563 [26] Y. Zheng, F. Liu, H.-P. Hsieh, U-Air: When urban air quality in-  
564 ference meets big data, 1436–1444, doi:\bibinfo{doi}{10.1145/2487575.  
565 2488188}, 2013.
- 566 [27] J. Kleine Deters, R. Zalakeviciute, M. Gonzalez, Y. Rybarczyk, Model-  
567 ing PM2.5 Urban Pollution Using Machine Learning and Selected Mete-  
568 orological Parameters, *Journal of Electrical and Computer Engineering*  
569 2017 (2017) 1–14, doi:\bibinfo{doi}{10.1155/2017/5106045}.
- 570 [28] Observatori del Turisme a Barcelona, Barcelona tourism activity re-  
571 port - 2018, [https://www.observatoriturisme.barcelona/sites/  
572 default/files/IAOTB18.pdf](https://www.observatoriturisme.barcelona/sites/default/files/IAOTB18.pdf), 2018.
- 573 [29] Eurostat Table [mar\_mg\_aa\_pwhd], Top 20 ports-gross weight of goods  
574 handled in each port, by direction, [https://appsso.eurostat.ec.  
575 europa.eu/nui/show.do?dataset=mar\\\_mg\\\_aa\\\_pwhd&lang=en](https://appsso.eurostat.ec.europa.eu/nui/show.do?dataset=mar\_mg\_aa\_pwhd&lang=en),  
576 2020.
- 577 [30] N. Perez, J. Pey, C. Reche, J. Cortés, A. Alastuey, X. Querol, Impact of  
578 harbour emissions on ambient PM10 and PM2.5 in Barcelona (Spain):  
579 Evidences of secondary aerosol formation within the urban area, *The  
580 Science of the total environment* 571 (2016) 237–250.

- 581 [31] J. Perdiguero, A. Sanz, Cruise activity and pollution: The case of  
582 Barcelona, *Transportation Research Part D: Transport and Environ-*  
583 *ment* 78 (2020) 102181, ISSN 1361-9209.
- 584 [32] G. Rodrguez, E. Martin-Alcalde, J. Murcia-González, S. Saurí, Eval-  
585 uating air emission inventories and indicators from cruise vessels at  
586 ports, *WMU Journal of Maritime Affairs* 16, doi:\bibinfo{doi}{10.1007/  
587 s13437-016-0122-8}.
- 588 [33] Departament de Territori i Sostenibilitat, Dades d'immissió  
589 dels punts de mesurament de la Xarxa de Vigilància  
590 i Previsió de la Contaminació Atmosfèrica, [https://analisi.transparenciacatalunya.cat/Medi-Ambient/  
591 Dades-d-immissi-dels-punts-de-mesurament-de-la-Xar/  
592 uy6k-2s8r](https://analisi.transparenciacatalunya.cat/Medi-Ambient/Dades-d-immissi-dels-punts-de-mesurament-de-la-Xarxa-uy6k-2s8r), portal de Dades Obertes de la Generalitat, 2020.
- 594 [34] Servei Meteorològic de Catalunya, Xarxa d'Estacions Meteorològiques  
595 Automàtiques (XEMA), [https://www.meteo.cat/observacions/  
596 llistat-xema](https://www.meteo.cat/observacions/llistat-xema), 2020.
- 597 [35] A. Tank, J. Wijngaard, G. Können, R. Böhm, G. Demarée, G. AA,  
598 M. Mileta, S. Pashiardis, L. Hejkrlik, C. Kern-Hansen, R. Heino,  
599 P. Bessemoulin, G. üller Westermeier, M. Tzanakou, S. Szalai,  
600 T. Pálsdóttir, D. Fitzgerald, S. Rubin, M. Capaldo, P. Petrovic, Daily  
601 surface air temperature and precipitation dataset 19011999 for Euro-  
602 pean Climate Assessment (ECA), *International Journal of Climatology*  
603 22 (2002) 1441 – 1453, doi:\bibinfo{doi}{10.1002/joc.773}.

- 604 [36] Barcelona Port Authority, Open Data of the Port of Barcelona, URL  
605 <http://opendata.portdebarcelona.cat>, 2020.
- 606 [37] Nicolas Bruno, Hourly flight movements at the Airport of Barcelona  
607 from 2015-2020, Private Communication between the authors and EU-  
608 ROCONTROL, 2020.
- 609 [38] Barcelona’s City Hall Open Data Service, Traffic state information by  
610 sections of the city of Barcelona, URL [https://opendata-ajuntament.  
611 barcelona.cat/data/en/dataset/trams](https://opendata-ajuntament.barcelona.cat/data/en/dataset/trams), 2020.
- 612 [39] B. B., M. Lang, L. Kotthoff, J. Schiffner, J. Richter, E. Studerus,  
613 G. Casalicchio, Z. M. Jones, mlr: Machine Learning in R, Journal of  
614 Machine Learning Research 17 (170) (2016) 1–5, URL [http://jmlr.  
615 org/papers/v17/15-066.html](http://jmlr.org/papers/v17/15-066.html).
- 616 [40] J. Benavides, M. Snyder, M. Guevara, A. Soret, C. Pérez García-Pando,  
617 F. Amato, X. Querol, O. Jorba, CALIOPE-Urban v1.0: Coupling R-  
618 LINE with a mesoscale air quality modelling system for urban air quality  
619 forecasts over Barcelona city (Spain), Geoscientific Model Development  
620 Discussions (2019) 1–35doi:\bibinfo{doi}{10.5194/gmd-2019-48}.
- 621 [41] J. Baldasano, M. Pay, O. Jorba, S. Gassó, P. Jimenez-Guerrero, An  
622 annual assessment of air quality with the CALIOPE modeling system  
623 over Spain, The Science of the total environment 409 (2011) 2163–78,  
624 doi:\bibinfo{doi}{10.1016/j.scitotenv.2011.01.041}.
- 625 [42] M. Guevara, F. Martínez, G. Arévalo, S. Gassó, J. Baldasano, An im-  
626 proved system for modelling Spanish emissions: HERMESv2.0, Atmo-

- 627 spheric Environment 81 (2013) 209 – 221, doi:\bibinfo{doi}{10.1016/j.  
628 atmosenv.2013.08.053}.
- 629 [43] W. Skamarock, J. Klemp, A time-split nonhydrostatic atmospheric  
630 model for research and NWP applications, J.Comput. Phys. 135.
- 631 [44] D. Byun, K. Schere, Review of the Governing Equations, Computa-  
632 tional Algorithms, and Other Components of the Models-3 Community  
633 Multiscale Air Quality (CMAQ) Modeling System, Applied Mechanics  
634 Reviews 59 (2006) 51–77.
- 635 [45] M. G. Snyder, A. Venkatram, D. K. Heist, S. G. Perry, W. B. Petersen,  
636 V. Isakov, RLINE: A line source dispersion model for near-surface re-  
637 leases, Atmospheric Environment 77 (2013) 748 – 756.

## 638 **Appendix A. Metrics**

The Geometric Mean Bias (GeoMean), the Fraction of Model Results Within a Factor of 2 of Observations (FAC2), the Correlation Coefficient ( $R$ ), the Mean Bias (MB) and the Root Mean Square Error (RMSE) are

defined as follows:

$$\text{GeoMean} = \exp \left( \overline{\ln \hat{Y}_i} - \overline{\ln Y_i} \right) \quad (\text{A.1})$$

$$\text{FAC2} = \sum_i^n \left[ 0.5 \leq \frac{\hat{Y}_i}{Y_i} \leq 2.0 \right] / n \quad (\text{A.2})$$

$$R = \frac{\overline{(\hat{Y}_i - \overline{\hat{Y}_i}) (Y_i - \overline{Y_i})}}{\sigma(Y_i) \sigma(\hat{Y}_i)} \quad (\text{A.3})$$

$$\text{MB} = \overline{Y_i - \hat{Y}_i} \quad (\text{A.4})$$

$$\text{RMSE} = \left( \overline{(Y_i - \hat{Y}_i)^2} \right)^{1/2} \quad (\text{A.5})$$

639 where  $Y_i$  and  $\hat{Y}_i$  are predicted and observed concentration data points. Aver-  
 640 age and standard deviation are respectively represented by the symbols  $\overline{(\dots)}$   
 641 and  $\sigma(\dots)$  and  $[\dots]$  are the Iverson brackets.

## 642 **Appendix B. Air and maritime traffic**

643 The distribution of number of simultaneously docked total vessels and  
 644 cruise ships at the Port of Barcelona over the October 2017 to March 2020  
 645 period (both included) are shown in Fig. B.7 and B.8 respectively.

646 Analogously, the distribution of number of flight operations at the Airport  
 647 of Barcelona over the October 2017 to March 2020 period (both included) is  
 648 shown in Fig. B.9.

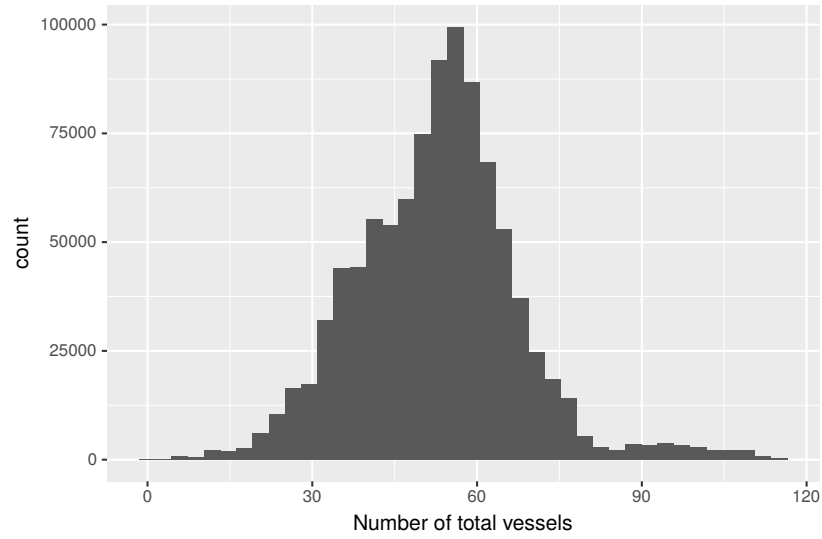


Figure B.7: Distribution of simultaneously docked total vessels at the Port of Barcelona from October 2017 to March 2020 both included.

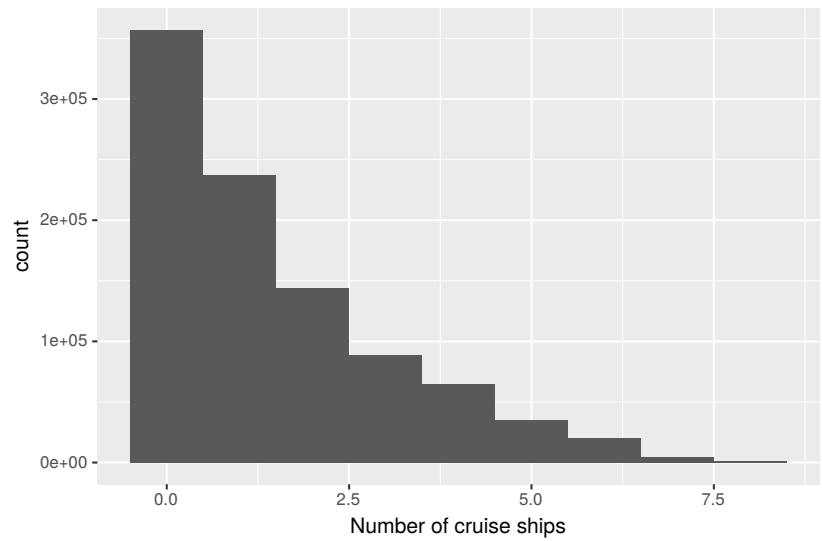


Figure B.8: Distribution of simultaneously docked cruise ships at the Port of Barcelona from October 2017 to March 2020 both included.

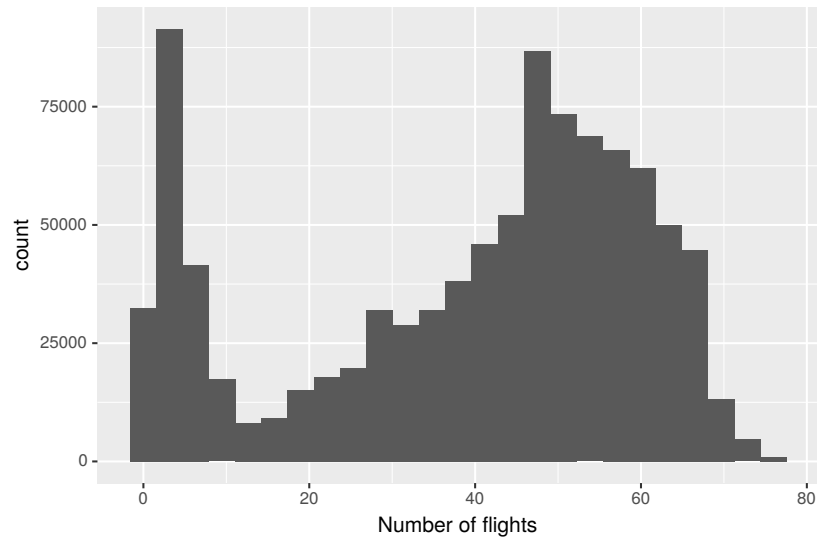


Figure B.9: Distribution of simultaneously docked cruise ships at the Port of Barcelona from October 2017 to March 2020 both included.

649 **Appendix C. Learner benchmark**

650 The comparison of performance for each optimized learner is shown in  
651 Fig. C.10 using the mean squared error for each pollutant averaged across  
652 stations.

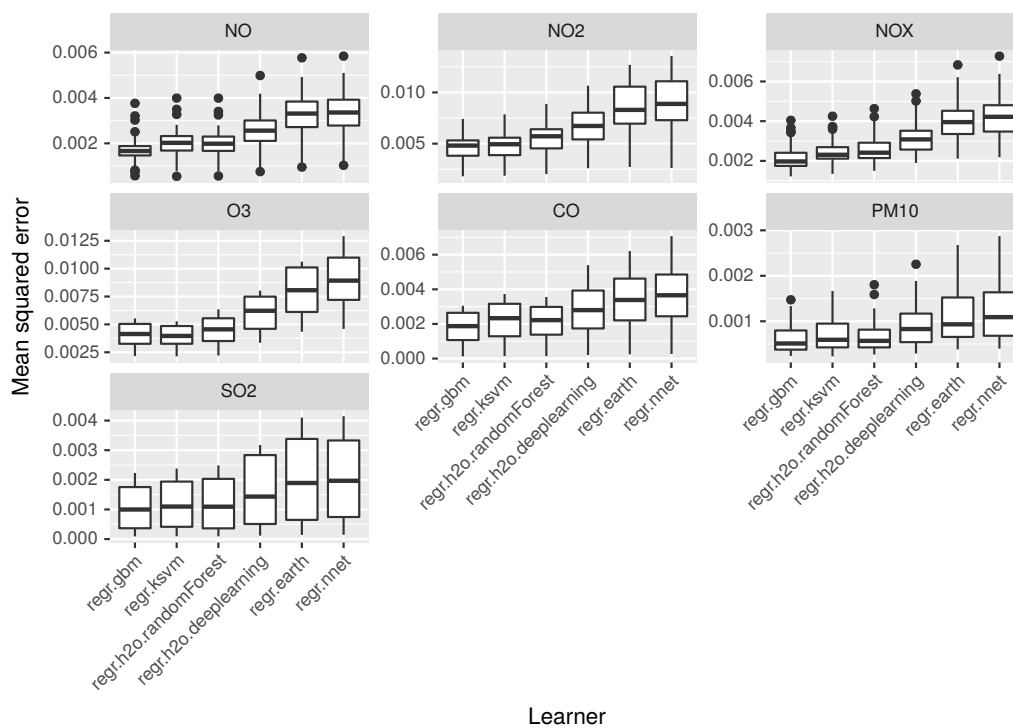


Figure C.10: Mean squared error for each pollutant across all stations for each learner using a 5-fold cross-validation resampling.