

### Abstract

$\text{NO}_2$  is a traffic-related air pollutant that is strongly associated with cardiovascular and respiratory diseases. Ground  $\text{NO}_2$  monitoring stations measure  $\text{NO}_2$  concentrations at certain locations and statistical predictive methods have been developed to predict  $\text{NO}_2$  as a continuous surface to inform decision-making. Among them, machine learning methods are the most powerful in capturing non-linear relationships between  $\text{NO}_2$  measurements and geospatial predictors, but it is unclear if the spatial structure of  $\text{NO}_2$  is also captured in the response-covariates relationships. In addition, most model comparison studies only compare accuracy in the prediction mean at ground stations, but do not consider prediction intervals and model interpretation and the effects of different model evaluation methods. In this study, we dive into the comparison between spatial and non-spatial data models accounting for the above-mentioned aspects. Moreover, we implemented a spatial and a non-spatial methods that have not been applied to air pollution mapping before and evaluated stack learning methods with and without modelling the spatial process. We implemented our study using national ground station measurements of  $\text{NO}_2$  in Germany and Netherlands of the year 2017, predicting  $\text{NO}_2$  to 100 m resolution grid. Our results indicate the importance of modelling the spatial process especially in areas close to traffic. The prediction intervals predicted with ensemble tree-based methods are satisfactory but too narrow with the geostatistical methods. Compared to ensemble tree-based methods, the geostatistical methods provide important spatial information for analysing emission sources and the spatial process of observations.

**Keywords:** geostatistics; machine learning; spatial prediction; model comparison; prediction interval; model interpretation

## **23    1    Introduction**

**24** NO<sub>2</sub> is a traffic-related air pollutant and has been found in epidemiological time series analysis  
**25** to highly associated with respiratory (Luo et al., 2016) and cardiovascular (Chiusolo et al., 2011)  
**26** diseases. NO<sub>2</sub> values are measured using monitoring stations at certain locations (e.g. close to traffic)  
**27** and most of the epidemiological studies identified the relationships between NO<sub>2</sub> and diseases or  
**28** hospital admission using a single NO<sub>2</sub> monitoring station to represent the entire district. However,  
**29** NO<sub>2</sub> is highly dynamic over the district and the difference in NO<sub>2</sub> concentrations will reflect on  
**30** personal exposures to NO<sub>2</sub>. Detailed spatial mapping of NO<sub>2</sub> is therefore required for more accurate  
**31** quantification of the relationships between NO<sub>2</sub> and health effects. In addition, detailed NO<sub>2</sub> maps  
**32** are necessary for scientific recommendations to be provided to policymakers and city planners.

**33** Statistical methods for NO<sub>2</sub> mapping have attracted a lot of attention with the burgeoning Ma-  
**34** chine Learning (ML)<sup>1</sup> methods and availability of ground monitoring station networks, atmospheric  
**35** satellite products, and geospatial predictors. Geospatial predictors are variables that are included  
**36** as covariates in a statistical air pollution model. Commonly used geospatial predictors are air  
**37** emission- (e.g. road networks) and dispersion-related (e.g. wind speed) variables, numerical mod-  
**38** elling (e.g. with chemistry transport model) output, and atmospheric remote sensing measurements  
**39** or products. A most recent (data available from Jan-2018) atmosphere sensing instrument, Tropomi

<sup>1</sup>list of abbreviations: CRPS: Continuous Ranked Probability Score; CV: Cross Validation; DF: Distributional Forest; GRF: Gaussian Random Field; GMRF: Gaussian Markov Random Field; GAMLS: Generalised Additive Models for Location Scale and Shape; INLA: Integrated Nested Laplace Approximation; IQR: Interquartile range; GWR: Geographic Weighted Regression; KED: Kriging with external drift; LUR: Land Use Regression; MAE: Mean Absolute Error; ML: Machine Learning; RF: Random Forest; OMI: Ozone Monitoring Instrument; Quantile Random Forest; RMSE: Root Mean Squared Error; SE: stacked ensemble; SPDE: Stochastic Partial Differential Equations; Tropomi: Tropospheric monitoring instrument; UK: Universal Kriging (UK); OMI (Ozone Monitoring Instrument) VIIRS: Visible Infrared Imaging Radiometer Suite; XGB: XGBoost

40 (Tropospheric monitoring instrument, NSO and ESA, 2019) onboard of Sentinel 5p satellite, mea-  
41 sures column density of a variety of gaseous air pollutants, in particular with an unprecedentedly  
42 high resolution for NO<sub>2</sub> (3.5 km x 5.5 km, across along track, since 06 August 2019).

43 Statistical methods applied for spatial air pollution prediction can be broadly classified depending  
44 on whether the spatial dependency is explicitly modelled. If not modelled, we refer to the methods  
45 "non-spatial" and otherwise "spatial". Most of the spatial air pollution models were developed to  
46 predict at coarser resolutions, commonly 1 km or coarser (Young et al., 2016; Shaddick et al., 2018;  
47 Beloconi and Vounatsou, 2020). Non-spatial methods are more dominant in air pollution mapping,  
48 particularly in high-resolution (100 m resolution or higher) mapping. Among them, LUR (Land  
49 Use Regression) models which assumes linear relationships between NO<sub>2</sub> and geospatial predictors  
50 are the most studied (Briggs et al., 2000; Hoek et al., 2008). Most recently, statistical learning (in  
51 this study, "statistical learning" is used interchangeably with "machine learning") methods (Hastie  
52 et al., 2009), including regularised linear regression (e.g. Lasso and Ridge regression (James et al.,  
53 2013)), kernel methods such as support vector machine (Suykens and Vandewalle, 1999), ensemble  
54 tree-based methods such as random forest (RF, Breiman, 2001) and XGBoost (XGB, Chen and  
55 Guestrin, 2016), have been applied for feature selection or capturing non-linear response-covariate  
56 relationships (Lu et al., 2020a; Chen et al., 2019a). In air pollution (not restricted to NO<sub>2</sub>) mapping,  
57 several studies compared between statistical learning and conventional LUR methods (Chen et al.,  
58 2019a; Kerckhoffs et al., 2019; Lu et al., 2020a; Ren et al., 2020; Rybarczyk and Zalakeviciute, 2018).

59 Geostatistical models (e.g. Kriging) and Geographically Weighted Regression (GWR) are the  
60 most used spatial methods for air pollution prediction (Vicedo-Cabrera et al., 2013; Li et al., 2014;  
61 Wang et al., 2021; Zou et al., 2016) and these methods have been combined with dimension reduction  
62 Zhai et al. (2018) and RF (Zhan et al., 2018; Liu et al., 2020) to improve NO<sub>2</sub> prediction accuracy.  
63 A Bayesian geostatistical model is developed in Beloconi and Vounatsou (2020) to predict NO<sub>2</sub> by

<sup>64</sup> integrating Tropomi satellite instrument measurements and chemical transport models. A GWR  
<sup>65</sup> model naturally models spatial varying coefficients by fitting multiple local regressions depending  
<sup>66</sup> on the homogeneity in response-covariate relationships when a number of observations are involved.  
<sup>67</sup> A typical geostatistical model can be viewed as consisting of two components: a mean function,  
<sup>68</sup> commonly a linear model, capturing the response-covariate relationships and a covariance function  
<sup>69</sup> modelling dependency of residuals from the mean (Bhatt et al., 2017). Conventional Kriging methods  
<sup>70</sup> suffer from the "big n problem", i.e. it may become computationally intractable with a large number  
<sup>71</sup> of observations. To deal with this problem, Lindgren et al. (2011) propose to use Stochastic Partial  
<sup>72</sup> Differential Equations (SPDE) to approximate the Gaussian Random Field (GRF) to a Gaussian  
<sup>73</sup> Markov Random Field (GMRF, Rue and Held (2005)). The main advantage of this is that the GMRF  
<sup>74</sup> has a sparse structure of the precision matrix, which is the inverse of the covariance matrix of a  
<sup>75</sup> GRF. Along with this, Rue et al. (2009) propose to use the Integrated Nested Laplace Approximation  
<sup>76</sup> (INLA) in a Bayesian framework to achieve the computational scalability of a geostatistical model  
<sup>77</sup> using approximations for all the estimations. This is especially advantageous when modelling NO<sub>2</sub>  
<sup>78</sup> over a larger scale e.g., continental or global-scale modelling when a large amount of observations  
<sup>79</sup> are modelled, and in spatiotemporal modelling.

<sup>80</sup> As spatial models are typically more complex compared to their non-spatial counterparts, several  
<sup>81</sup> studies compared spatial and non-spatial models to understand if the spatial effects could be simply  
<sup>82</sup> modelled by including certain covariates in LUR models. Young et al. (2016) studied the use of  
<sup>83</sup> universal Kriging (UK), OMI (Ozone Monitoring Instrument) satellite instrument (Earthdata) and  
<sup>84</sup> LUR models for NO<sub>2</sub> prediction at 2.5 km resolution. Young et al. (2016) indicated that either  
<sup>85</sup> using UK or adding OMI in the LUR model improves a LUR model but adding OMI in a UK  
<sup>86</sup> model only trivially improves the performance. Bertazzon et al. (2015) shows that the inclusion of  
<sup>87</sup> the meteorological variables accounts for spatial effects similarly to the use of spatial autoregressive

88 models(Anselin et al., 2001). However, even if the spatial dependency can be captured by involving  
89 certain covariates in a LUR model, we may still need geostatistical methods to understand the  
90 spatial structure present in the data. Linear models have been used for the mean function but the  
91 relationships between NO<sub>2</sub> and predictors have been shown to be better modelled with non-linear  
92 ML methods (Lu et al., 2020a). Most recent studies attempt to replace the linear mean function  
93 with ML models. Liu et al. (2020) applied a geostatistical model to the residuals from an RF model  
94 for the spatial prediction of PM<sub>2.5</sub>. In disease mapping, Bhatt et al. (2017) proposes to stack ML  
95 models to replace the mean function in a geostatistical model.

96 Few studies have compared between geostatistical and ML methods, possibly because the ML  
97 methods are still relatively less studied in air pollution mapping and in the field of geostatistics. It  
98 might be more interesting to compare between geostatistical methods and ML methods than geosta-  
99 tistical methods and LUR, because ML methods may be more capable of (implicitly) capturing the  
100 spatial dependency by integrating covariates, when the number of observations is sufficient. More-  
101 over, most comparison studies only compare the cross-validation accuracy of the prediction mean  
102 (e.g. using R-squared, mean absolute error, or root mean squared error), ignoring the prediction  
103 intervals. Also not discussed is the cause of the prediction errors, are they caused by missing co-  
104 variants, violation of the model assumptions (e.g. data distribution, non-linearity), or inconsistent  
105 distributions between training and validation sets. Also, different cross-validation strategies, e.g.,  
106 how do we split the train-test sets, may lead to different model validation results. Current studies  
107 typically solely rely on k-fold splitting (Kerckhoff et al., 2019; Larkin et al., 2017; Ren et al., 2020)  
108 or bootstrapping (Lu et al., 2020a) to randomly splitting between train-test sets, which may be  
109 one-sided and does not provide an indication of accuracy in spatial blocks (but only at the locations  
110 of ground stations).

111 In this study, we focus on ensemble tree-based methods (e.g. RF and boosting) in the ML

category and a hierarchical spatial model (Lindgren et al., 2015; Blangiardo and Cameletti, 2015; Moraga, 2019) called latent Gaussian model in the geostatistics category. Additionally, we invest in stacked models in integrating ML and geostatistical models and develop a LUR model using Lasso for comparison. Ensemble trees are nonparametric models, deriving prediction intervals is therefore less straightforward than a parametric model (e.g. a linear regression model) but has been studied and shown satisfactory results with simulated data. Prediction intervals have been most well studied for RF (Meinshausen, 2006; Wager et al., 2014; Stasinopoulos et al., 2007; Alakus et al., 2021) and more recently for boosting (Duan et al., 2020; Velthoen et al., 2021). Comparing probabilistic methods (i.e. prediction interval calculation) of RF and boosting is beyond the scope of this study and we focus on prediction intervals derived for RF to compare with geostatistical methods. Possibly, one of the most widely recognisable methods to derive RF prediction intervals is Quantile Random Forest (QRF) (Meinshausen, 2006). QRF has been shown to estimate middle quantiles well but may fall short at the extremes due to the limited number of observations in the tail regions (Velthoen et al., 2021). Velthoen et al. (2021) proposed to use extreme quantile regression to estimate for data outside the range of observations. Another well-recognised method is distributional regression forests (DF) (Schlosser et al., 2019), which embeds the GAMLSS (Generalised Additive Models for Location Scale and Shape) (Stasinopoulos et al., 2007) into RF.

Fouedjio and Klump (2019) compared prediction accuracy and uncertainty quantification between KED (Kriging with external drift) and QRF by simulating data with various levels of spatial dependency. It concluded that an optimal model choice depends on the level of spatial dependency and response-covariate relationships. However, it does not account for the fact that in practice, as an ensemble tree-based method can make use of a large number of (possibly correlated) predictors without being constrained to certain (e.g. linear) relationships, the spatial dependency may be explained by the covariates despite not being explicitly modelled.

<sup>136</sup> The objective of our study is to compare geostatistics and non-spatial ensemble tree-based models  
<sup>137</sup> for NO<sub>2</sub> mapping, in terms of their prediction accuracy, uncertainty quantification, and model inter-  
<sup>138</sup> pretation and to understand effect of modelling spatial structures. More specifically, the following  
<sup>139</sup> sub-objectives are reached:

- <sup>140</sup> 1. Optimising a set of spatial hierarchical and ML models for NO<sub>2</sub> prediction in Germany and  
<sup>141</sup> the Netherlands.
- <sup>142</sup> 2. Developing a non-spatial and a geostatistical stacked ensemble model, i.e., a stack of various  
<sup>143</sup> ML learners.
- <sup>144</sup> 3. Model comparison regarding the predicted mean, prediction interval, and model interpretation.

<sup>145</sup> The spatial Hierarchical model incorporates the spatial random effect along with other covariates  
<sup>146</sup> and the estimation is performed using the R package **INLA** (Rue et al., 2009; Martins et al., 2013).  
<sup>147</sup> XGB, RF and Lasso are chosen for the comparison with the geostatistical model and they also  
<sup>148</sup> form the base learners in the two (geostatistical and non-spatial) stacked learning models. The ML  
<sup>149</sup> methods are chosen for their dissimilarity. Specifically, Lasso is a linear regression model without  
<sup>150</sup> accounting for spatial dependency. RF and XGB are non-linear models with regression trees as base-  
<sup>151</sup> learners and are not affected by dependent covariates. XGB is a highly scalable boosting method  
<sup>152</sup> that builds tree models subsequently over the residuals of previous trees and has multiple routines  
<sup>153</sup> to penalise model over-fitting (Chen et al., 2019b), which has been reported in various studies to  
<sup>154</sup> obtain the highest prediction accuracy Lu et al. (2020a).

## <sup>155</sup> 2 Data

<sup>156</sup> NO<sub>2</sub> concentration measurements of 2017 from national ground stations of Germany and the Nether-  
<sup>157</sup> lands are used. The original hourly data is downloaded from the EEA (European Environment

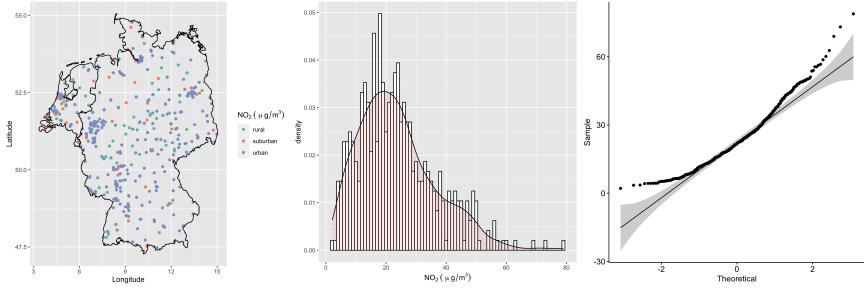


Figure 1: Spatial distribution of NO<sub>2</sub> stations, histogram and Q-Q plot of the NO<sub>2</sub> measurements.

158 Agency, Nelson, 1999; EEA, 2021). Negative values are considered as missing. The data is ag-  
 159 gregated to annual concentrations by taking the mean and omitting missing values. The spatial  
 160 distribution of NO<sub>2</sub> stations and the station types, histogram and Q-Q plot for normality are shown  
 161 in fig. 1. We conducted a Shapiro test for normality, with the result implying the distribution of data  
 162 being significant different from normal distribution ( $p\text{-value}= 8.605\text{e-}12$ , "normal distribution" and  
 163 "Gaussian distribution" are used interchangeably in this study). A Gamma distribution test was  
 164 conducted using the method proposed in Villaseñor and González-Estrada (2015) and implemented  
 165 in Gonzalez-Estrada and Villasenor-Alva (2020). The test result ( $p\text{-value} = 0.32$ ) indicates that the  
 166 data distribution is not significantly different from Gamma distribution.

167 The geospatial predictor grids (table 1) are calculated or re-sampled at 100 m resolution. They  
 168 are either spatial attributes aggregated in a circular ring centred at each sensor or prediction location,  
 169 called buffered predictors, or values of the spatial attribute at the observation or prediction location,  
 170 called gridded variables. The buffered predictors include total road length, total industry areas,  
 171 VIIRS (Visible Infrared Imaging Radiometer Suite) Nighttime Day/Night Band radiances values  
 172 (nightlight, NOAA, 2021) and population. Variables that are originally grids include wind speed  
 173 and temperature (Dee et al., 2011), elevation (NASA), annual mean Tropomi level 3 product of NO<sub>2</sub>  
 174 column density (Copernicus, 2021) from 2019 (due to the increased resolution compared to 2018).

<sub>175</sub> The buffered predictors of road and industry are calculated from OpenStreetMap (OpenStreetMap  
<sub>176</sub> contributors, 2019). For detailed descriptions of the processing of the geospatial predictors please  
<sub>177</sub> refer to Lu et al. (2020a).

### <sub>178</sub> 3 Methods

<sub>179</sub> The methods considered in this study are classified as spatial and non-spatial and are given the  
<sub>180</sub> names below in this study.

#### <sub>181</sub> **Spatial models:**

- <sub>182</sub> 1. INLA: A spatial hierarchical model fit using INLA with a Gaussian likelihood.
- <sub>183</sub> 2. INLA-G: A spatial hierarchical model fit using INLA with a Gamma likelihood.
- <sub>184</sub> 3. SE-INLA: using the spatial hierarchical model to stacked learning with Lasso, RF and XGB  
<sub>185</sub> models as base learners;

#### <sub>186</sub> **Non-spatial models:**

- <sub>187</sub> 1. LA: A Lasso regression model;
- <sub>188</sub> 2. RF: A RF model;
- <sub>189</sub> 3. XGB: An XGB model assuming a Gaussian objective function;
- <sub>190</sub> 4. XGB-G: An XGB model assuming a Gamma objective function;
- <sub>191</sub> 5. QRFLA: using Lasso to aggregate QRF trees (Hastie et al., 2009);
- <sub>192</sub> 6. SE: stacked learning with Lasso, RF and XGB models as base learners;
- <sub>193</sub> 7. QRF: quantile regression forest (Meinshausen, 2006);
- <sub>194</sub> 8. DF: distributional regression forest (Schlosser et al., 2019).

Table 1: Geospatial predictors considered in this study. ”\_mon” indicates months (mon = 1, 2....,12). ”\_buf” indicates buffer radius in meters. The road length and industrial areas are calculated with buffer radii of 100 m, 300 m, 500 m, 800 m, 1000 m, 3000 m and 5000 m. The night lights digital numbers are calculated with buffer radii of 450 m, 900 m, 3150 m and 4950 m. The original resolution is provided for gridded variables and data types for vector variables.

Predictor	Variable name	Unit	Resolution/data type
Monthly wind speed at 10 m altitude.	Wind_speed_10m_mon	km/hr	10 km
Monthly temperature at 2 m altitude.	temperature_2m_mon	Celsius	10 km
TROPOMI 2018 mean vertical column density.	trop_mean_filt; Tropomi	$mol/cm^2$	0.01 arc degrees
Population in 5 km grid	population_5000	count	5 km
Population in 3 km grid	population_3000	count	3 km
Population in 1 km grid	population_1000	count	1 km
Nightlight	nightlight_bufnl	$Wcm^{-2}sr^{-1}$	500 m
Total length of highway	road_1_buf	m	polygon, lineString
Total length of primary roads	road_2_buf	m	polygon, lineString
Total length of local roads	road_M345_buf	m	polygon, lineString
Area of industry	I_1_buf	$m^2$	polygon, lineString

195 **3.1 Non-spatial methods**

196 Lasso is a linear regression algorithm with the L1 regularisation to shrink variable coefficients to  
197 zero, which enables "feature selection". In the cost function, the absolute value of coefficient is added  
198 to the original least squares as a penalty term. RF and XGB in this study use trees as base learners  
199 and ensemble them to reduce variability of single trees (Friedman, 2001). RF firstly randomly draws  
200 a subset of features, and then choose features from this subset to build the tree. RF (Breiman, 2001)  
201 grows trees independently and then take the mean of the predictions of each tree.

202 QRF is a non-parametric prediction interval estimation method which keeps all the observations  
203 in the terminal node for estimating the conditional probability function. Specifically, it samples  
204 from all the response values in each terminal node and use the ratio between the number of samples  
205 that is taken from each terminal node and the number of total observations in the terminal node as  
206 weights to aggregate the samples. The weights of all the trees are summed. The summed weights  
207 computed for each observation are then used to construct the empirical conditional cumulative  
208 distribution function (Meinshausen, 2006). QRFLA uses Lasso as a post-processing of QRF (Hastie  
209 et al., 2017, page 617). This method firstly preserves all the trees instead of aggregating them  
210 (e.g. taking the mean of all the predictions) and then apply Lasso regression to all the trees for  
211 aggregation. This leads to a shrinkage of the tree space and theoretically reduces model variance.  
212 DF (Schlosser et al., 2019) firstly divide data into regions as homogeneous as possible with respect  
213 to a parametric distribution, thus capturing changes in location, scale and shapes. For each tree,  
214 maximum likelihood is used to fit distributions and recursively select and split covariates according  
215 to the instability of the gradient of the likelihood at each observation along each co-variate. Then,  
216 the distributional trees are ensembled for DF.

217 XGB is a variation of gradient boosting, which grows trees subsequently by fitting to model  
218 residuals of the previous step. XGB is scalable to multiple threads. It enables multiple penalisation

219 paths to control model complexity to prevent model over-fitting, including regularisation (e.g. L1  
220 regularisation) on tree width and terminal node values, as well as drop-out (dropping trees), sampling  
221 observations (take a subset of observations in each run), and early stopping (stop iterating when after  
222 a few rounds the loss does not decrease or the node does not meet the splitting rule). The default  
223 objective function for regression assumes normal distribution of target variables (and the prediction  
224 is the mean of the distribution). This assumption is used in all the air pollution mapping studies.  
225 Here, we additionally fit a model with the objective function assuming the target variable follows a  
226 Gamma distribution (XGB-G) as the distribution of NO<sub>2</sub> measurements is closer to Gamma than  
227 normal distribution.

228 Different from the ensembling in RF or XGB,SE (Stacking Ensemble) refers to a class of al-  
229 gorithms that trains a second-level “meta-learner” to optimise the combination of a collection of  
230 prediction algorithms (base-learners). The base-learners are preferably diverse to capture different  
231 relationships or patterns. In this study, Lasso, RF, and XGB are the base-learners. Cross-validated  
232 predicted values (commonly known as level-one data) are used to train the meta-learner.

### 233 **3.2 Hyperparameter setting for XGB and RF**

234 To optimise the hyperparameters of XGB (known as ”model tuning”), we used grid search to optimise  
235 hyperparameters in 5-fold cross-validation basing on the minimum RMSE (Root Mean Squared  
236 Error) and additionally manual adjustment of the hyperparameters to look at the prediction patterns.  
237 The grid search is used instead of more computationally efficient methods (e.g. Bayesian or random  
238 search) as the optimal hyperparameter range is largely known from our previous experiences (Lu  
239 et al., 2020a, 2021). The search grid for the number of iterations (nrounds) was from 200 to 3000,  
240 with a step of 200; maximum tree depth (max-depth) from 3 to 6 with a step of 1, learning rate  
241 (eta) from 0.001 to 0.1 with a step of 0.05, the penalty term Gamma (Chen et al., 2019b) from 1

242 to 5 with a step of 1, the subsample is set to 0.7, L1 norm penalisation (lambda) is set to 2 and L2  
 243 norm penalisation (alpha) is set to 0. RF is not sensitive to hyperparameter tuning. We used the  
 244 default setting of number of variables that are randomly drawn for each tree (Breiman, 2001), which  
 245 is the integer part of the total number of variables divided by three. The number of trees is set to  
 246 2000 for a safe choice as the high number of trees will not negatively affect model performance.

### 247 3.3 Geostatistical models

248 Suppose we assume that  $NO_2$  values  $y_i$  measured at locations  $\mathbf{s}_i$ ,  $i = 1, \dots, n$ , follows a Gaussian  
 249 distribution with mean  $\mu_i$  and variance  $\sigma^2$ , where the mean  $\mu_i$  is expressed as a sum of covariates  
 250 and a spatially structured random effect following a zero-mean Gaussian process with a spatial  
 251 covariance function (Moraga, 2019).

$$y_i \sim N(\mu_i, \sigma^2), \quad i = 1, 2, \dots, n \quad (1)$$

$$\mu_i = \mathbf{d}_i \boldsymbol{\beta} + \mathbf{x}(\mathbf{s}_i) \quad (2)$$

252 Here,  $\mathbf{d}_i = (d_{i1}, \dots, d_{ip})$  is the vector of covariates at location  $\mathbf{s}_i$ ,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$  is the  
 253 coefficient vector, and  $\mathbf{x}(\mathbf{s}_i)$  denotes a spatial Gaussian random field. That is,  $\{\mathbf{x}(\mathbf{s}_1), \dots, \mathbf{x}(\mathbf{s}_n)\} \sim$   
 254  $\mathcal{N}_n(\mathbf{0}, \boldsymbol{\Sigma})$ , where  $N_n$  is a Normal multivariate distribution for the spatial process specified by its  
 255 mean  $\mathbb{E}(\mathbf{x}(\mathbf{s}))$ , and covariance function  $C(\mathbf{s}_1, \mathbf{s}_2) = \text{Cov}(\mathbf{x}(\mathbf{s}_1), \mathbf{x}(\mathbf{s}_2))$ . The Gaussian random field  
 256 can be stationary and isotropic, where the covariance function depends only on the distance and not  
 257 direction between points, that is  $C(\mathbf{s}_1, \mathbf{s}_2) = \text{Cov}(\|\mathbf{s}_1 - \mathbf{s}_2\|)$  and its dependence is commonly modeled  
 258 using a Matérn function (Stein (2012); Yuan (2011); Diggle et al. (2013)). Since incorporating the  
 259 spatial dependence directly with a large number of observations using a Gaussian random field is  
 260 computationally expensive, Rue and Held (2005) proposed the approximation of a Gaussian random  
 261 field by a Gaussian Markov random field for a more efficient computational process of estimation.

262 The main property of the Gaussian Markov random field is that it uses a conditional dependency  
263 structure through the precision matrix  $\mathbf{Q}$ .

264 In this study, we compare two spatial hierarchical models with geospatial predictors as covariates,  
265 one uses a Gaussian likelihood and the other a Gamma likelihood. The Gamma model has the same  
266 hierarchical structure as the Gaussian model: the response variable in (1) can be represented by  
267  $y_i \sim \text{Gamma}(\alpha, \beta)$  where  $\alpha$  is the shape parameter and  $\beta$  the rate parameter. The INLA-SE model  
268 uses a Gaussian likelihood.

269 **3.4 INLA and SPDE**

270 To fit the geostatistical models, we use the R package **INLA** which facilitates the application of the  
271 INLA and the SPDE approaches. Following the expression proposed in (1), the structure for the  
272 hierarchical model is:

$$\mathbf{y} | \mathbf{x}, \theta_1 \sim N(\mathbf{D}\boldsymbol{\beta} + \mathbf{A}\mathbf{x}, \theta_1) \quad (3)$$

$$\mathbf{x} | \theta_2 \sim \text{GRF}(\mathbf{0}, \mathbf{Q}(\theta_2)^{-1}) \quad (4)$$

$$\boldsymbol{\theta} = \{\theta_1, \theta_2\} \quad (5)$$

273 where  $\boldsymbol{\theta}$  is the vector of hyperparameters with  $\theta_1 = \sigma^2$ ,  $\theta_2 = \{\log(\tau), \log(\kappa)\}$ ,  $\mathbf{x}$  is the spatial  
274 latent field,  $\mathbf{A}$  is the projector matrix and  $\mathbf{y}$  is the vector of the response variable  $f(\cdot | \mathbf{x}, \boldsymbol{\theta})$ ,  
275 commonly from the exponential family of distributions.  $\mathbf{D}$  is a covariate matrix and  $\boldsymbol{\beta}$  a coefficient  
276 matrix.

277 The R package **INLA** can be used to perform direct numerical calculation of the posterior distri-  
278 bution for a Bayesian hierarchical model (Rue et al. (2009), Martino and Rue (2009)). If we use  $\mathbf{x}$   
279 as a latent Gaussian field (a Gaussian Markov random field),  $\boldsymbol{\theta}$  a vector of (hyper)parameters and

280  $\mathbf{y}$  a vector of observations, assuming independent observations given the vector of the spatial latent  
281 field ( $\mathbf{x}$ ) and the hyperparameters ( $\boldsymbol{\theta}$ ), the likelihood can be expressed as:

$$p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) = \prod_{i \in \mathcal{I}} p(y_i \mid \eta_i, \boldsymbol{\theta}), \quad (6)$$

282 where  $\eta_i$  is the linear predictor and  $\mathcal{I}$  contains the indices of the observed values  $\mathbf{y}$ .

283

284 The main aim is to approximate the posterior density for the posterior of the spatial latent field  
285 and the hyperparameters. The marginal densities can be obtained:

$$p(x_i \mid \mathbf{y}) = \int p(x_i \mid \boldsymbol{\theta}, \mathbf{y}) p(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}, \quad (7)$$

286 and

$$p(\boldsymbol{\theta}_j \mid \mathbf{y}) = \int p(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}_{-j}. \quad (8)$$

287 respectively (Lindgren et al. (2015); Krainski et al. (2018)).

288

289 To model data indexed in space, Lindgren et al. (2011) proposed a new methodology based mainly  
290 on the approximation of the Gaussian random field with the Matérn function using the Stochastic  
291 Partial Differential Equations (SPDE) as follows:

$$(\kappa^2 - \Delta)^{\alpha/2}(\tau(\mathbf{s})x(\mathbf{s})) = \mathbf{W}(\mathbf{s}), \quad (9)$$

292 where  $\kappa$  is a scale parameter,  $x(\mathbf{s})$  is a spatial random field,  $\Delta$  is the Laplacian,  $\alpha$  is the parameter  
293 that controls the smoothness of the realizations,  $\tau$  controls the variance and  $\mathbf{W}(\mathbf{s})$  is a Gaussian  
294 spatial white noise process (Lindgren et al. (2015)). For the above we can use a Gaussian Markov  
295 random field that approximates to a Gaussian random field using a triangulation of the region of

296 study without specifying an explicit covariance structure through the SPDE method. This approx-  
297 imation leads to a decrease in computational burden from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n^{3/2})$ .

### 298 3.5 Geospatial predictor selection for the INLA model

299 As involving too many covariates (e.g. more than 12) in the INLA model brings problems in model  
300 inferencing and multicollinearity, we used Lasso to reduce the number of variables. The Lasso was  
301 used instead of ensemble tree-based methods for feature selection because it is also a linear model  
302 (same as the INLA and INLA-G models in our study). Variables are selected with the L1 norm  
303 penalty that returns a model with errors that are within one standard error of the minimum mean  
304 cross-validated error. Lasso is applied to 80% data randomly sampled from all the observations  
305 and this process is repeated 20 times. Variables that are selected more than 90% of the times (i.e.  
306 18) will be considered as covariates in INLA. The times that the Lasso selected certain variables is  
307 shown in table 2. The INLA modelling process applies the same bootstrapped samples for training  
308 and validation. In addition, AIC (step-wise) model selection is applied to the entire dataset to  
309 suggest a model as a further reference. The variables selected by AIC are almost the same as  
310 Lasso selected variables, besides it does not choose road\_class\_3\_3000, which is highly correlated  
311 with road\_class\_1\_5000. Based on this, the road\_class\_3\_3000 is not used as a covariate in INLA.

### 312 3.6 INLA model parameterisation

313 The triangulated mesh constructed in the SPDE approach is shown in supplementary material  
314 (supfig. 1), with size of the inner and outer extensions around the data locations (*offsets*) 1/8 of  
315 the maximum distance among all the observations for both the inner and outer extensions. The  
316 maximum allowed triangle edge lengths in the region and in the extension (*max.edge* ) are set  
317 to respectively 1/30 and 1/5 times maximum distance among all the observations. The Matern

Table 2: Frequency (number of times) of variables selected by Lasso in 20 times bootstrapping and variables that are selected more than 90% times (i.e. 18) are listed below. These variables are considered in INLA besides road\_class\_3\_3000.

	Variables	Frequency
1	nightlight_450	20
2	population_1000	20
3	population_3000	20
4	road_class_1_5000	20
5	road_class_2_100	20
6	road_class_3_300	20
7	trop_mean_filt	20
8	road_class_3_3000	19
9	road_class_1_100	18

318 SPDE model is constructed with  $\alpha = 2$ . The SE-INLA model has the same specification (i.e.  
319 mesh structure, likelihood, objective function, priors, optimisation process) as the INLA model  
320 parameterisation described above.

## 321 4 Model evaluation

### 322 4.1 Cross validation

323 We use RMSE, MAE (Mean Absolute Error), IQR (Interquartile Range) and R<sup>2</sup> (R-squared) to  
324 compare model performance. RMSE is calculated as the square root of the differences between  
325 predictions and observations; MAE is calculated as the absolute differences between predictions  
326 and observations; IQR is the differences between the third and first quartiles of the prediction. R<sup>2</sup>  
327 indicates the explained variance and is calculated as  $R^2 = 1 - \text{var}(\text{error})/\text{var}(y)$ , where var(.)  
328 indicates variance, error indicates model residuals and y indicates observed response values. When  
329 different data is used in CV (e.g. separating between close and far-away from roads), we additionally  
330 calculated the RRMSE (relative RMSE), RMAE (relative MAE), RIQR (relative IQR) to account  
331 for the differences in the magnitudes of response values. The RRMSE and RMAE are calculated by  
332 dividing the RMSE and MAE, respectively, by the mean of observations. The RIQR was calculated  
333 by dividing the IQR by the median of observations. The three CV methods we designed and used  
334 to assess our model performance are:

- 335 1. Bootstrapped CV. 20-times randomly bootstrapped splitting of training and test sets (Lu et al.,  
336 2020a).
- 337 2. Spatial-blocked CV. Dividing data into spatial blocks, each time use one block for test and  
338 other blocks for training.

339     3. Customised CV. Splitting train-test based on values of certain covariates. In this study, three  
340       sub-areas are defined, 1) close to traffic and with high population ("tr-hp"), 2) close to traffic  
341       and with middle low population ("tr-lmp"), 3) far away from traffic ("far"). High population is  
342       defined as the variable population of 1000 m buffer that is in the last quartile. Low population  
343       is defined as the variable population of 1000 m buffer is below the median. Close to road is  
344       defined as (please refer to table 1 for the definition of covariates):

```
345       road_class_2_100 > 0 |  
346       road_class_1_100 > 0 |  
347       road_class_3_100 > quantile(road_class_3_100, .75))
```

348     Far away from road is defined as:

```
349       road_class_2_100 == 0 &  
350       road_class_1_100 == 0 &  
351       road_class_3_100 < quantile(road\class\_3\_100, .5)
```

352     where "&" indicates "and" and "|" indicates "or". The second variable of the function  
353     "quantile(.)" indicates the percentage quantile of the variables.

354     This yields 85, 65, and 177 samples in each category. This ensures a balanced number of samples  
355     between close to traffic and far-away from traffic. Each time, 30 samples (7% of the entire dataset)  
356     are drawn from the corresponding category for CV. For example, each time, 30 samples are drawn  
357     from the 85 samples as the test set to obtain the prediction accuracy CV for the situation "tr-hp"  
358     and the rest is used for training.

## 359     **4.2 Prediction intervals**

360     CRPS (Continuous Ranked Probability Score) and coverage probabilities are used as quality indica-  
361     tors of prediction intervals. CRPS is an uncertainty measure that assesses the similarities between

362 two distributions. We use it to indicate how the predicted distribution matches the observed dis-  
363 tribution. The CRPS implemented as an R package **ScoringRules** (Jordan et al., 2017) is used.  
364 CRPS is calculated for the INLA and QRF models. For the INLA model, the prediction intervals  
365 are calculated by simulating from the response  $Y \sim N(\theta, \sigma^2)$  where  $\theta$  and  $\sigma^2$  are the fitted mean and  
366 variance. The mean of CRPS for all the points within each test block is calculated in spatial-blocked  
367 CV. Coverage probabilities are calculated as the ratio between the number of predictions within  
368 the upper and lower quantile and the total number of predictions (in the test set). The prediction  
369 intervals are mainly compared between INLA, INLA-G, QRF and DF. The prediction interval for  
370 QRFLA is compared with QRF to investigate the effects of Lasso tree-aggregation strategy on the  
371 prediction intervals.

### 372 4.3 Model interpretation

373 We inspect fixed and spatial random effects modelled by INLA and compare the spatial random field  
374 with modelled prediction intervals and model residuals to understand the contribution of spatial  
375 random effects. Different from linear regression methods, which themselves are the best models for  
376 interpretation, interpreting ensembling tree-based methods requires external models (Lundberg and  
377 Lee, 2017). We use SHAP (SHapley Additive exPlanations, Lundberg et al., 2018; Lundberg and  
378 Lee, 2017), a unified method based on additive feature attribution, to estimate variable influence in  
379 RF and XGB models.

## 380 5 Results

### 381 5.1 Accuracy assessment and uncertainty quantification

#### 382 Non-spatial CV

383 Both ensemble tree-based methods with a Gaussian objective function and INLA with a Gaussian  
384 likelihood function obtain higher prediction accuracy than Lasso (table 3), indicating the necessity of  
385 using a more flexible model and modelling spatial random fields. Among individual methods, in terms  
386 of  $R^2$  and RMSE, INLA with Gaussian likelihood obtained the highest prediction accuracy, followed  
387 by XGB-G and QRFLA. QRFLA greatly improves from original RF. Despite the distribution of  
388 response being closer to Gamma distribution compared to Gaussian distribution, using Gamma  
389 regression in XGB and specifying Gamma likelihood in INLA both decrease the prediction accuracy  
390 considerably. Compared to INLA, XGB obtained lower RMSE and  $R^2$  despite it obtained lower  
391 MAE and IQR, indicating that the XGB model predicts less well at more extreme ranges. The  
392 QRF and DF results are not shown in table 3 as the results are very similar to RF. Their prediction  
393 intervals are compared.

394 SE-INLA improves prediction accuracy compared to SE and INLA, obtained the best results in  
395 terms of root mean squared error (6.83, 24.5% of the mean of observations) and  $R^2$  (0.71). This in-  
396 dicates the spatial structures could further improve prediction accuracy despite flexible relationships  
397 captured from ML models.

### 398 **Spatial-blocked CV**

399 Spatial-blocked CV provides information about prediction accuracy in spatial blocks. The  $R^2$   
400 map (fig. 2) shows that the XGB, RF and INLA predict relatively well in most parts of Germany  
401 besides blocks at the boundaries. The  $R^2$  for the block western the Netherlands is also relatively low  
402 with all the three methods and especially for XGB ( $R^2$ : 0.2). RF obtains the best result for the block  
403 of western the Netherlands ( $R^2$ : 0.5). The INLA model outperforms RF and XGB in the blocks  
404 at south-east and north. The  $R^2$  between blocks are the most heterogeneous with XGB, which is  
405 consistent to the result of bootstrapped CV that the XGB falls short at predicting extremes.

406 The spatial-blocked CRPS fig. 3 is computed for QRF and INLA (the DF is not shown as it will

Table 3: Prediction accuracy matrix for different models using 20 times bootstrapped cross-validation. Non-spatial models: LA: Lasso; RF: random forest, XGB: XGBoost using the default Gaussian loss; XGB-G: XGBoost using a Gamma loss; QRFLA: quantile random forest with Lasso for shrinkage aggregation of regression trees; SE: stacked ensembling. Spatial models: INLA: a latent Gaussian model implemented using INLA assuming a Gaussian likelihood. INLA-G: a latent Gaussian model implemented using INLA assuming a Gamma likelihood. SE-INLA, geostatistical stacked ensembling.

	LA	RF	XGB	XGB-G	QRFLA	SE	INLA	INLA-G	SE-INLA
RMSE	7.54	7.45	7.14	8.91	7.23	7.18	7.06	9.21	6.83
IQR	8.47	7.39	6.54	9.21	7.27	7.30	7.1	7.4	6.8
MAE	5.69	5.51	5.05	6.27	5.28	5.31	5.3	6.2	5.0
R <sup>2</sup>	0.65	0.65	0.68	0.51	0.67	0.69	0.69	0.45	0.71

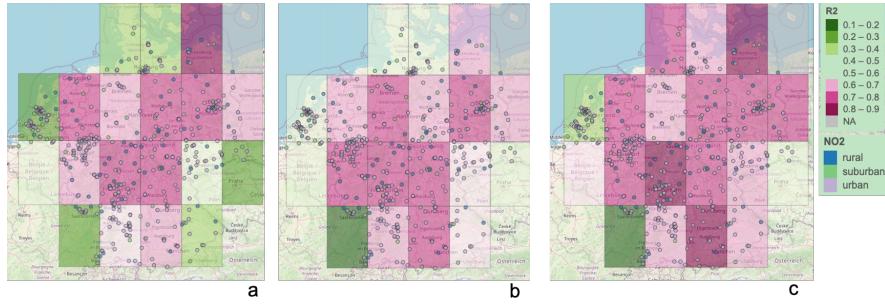


Figure 2: The R-squared of each block, using the rest of the blocks for training. The models are a) XGB, b) QRF, c) INLA.



Figure 3: The CRPS (Continuous Ranked Probability Score) of each block, using the rest of the blocks for training. a) RF, b) INLA.

be shown that the QRF and DF performed similarly in prediction interval prediction (section 5.2)).  
 The INLA predicted prediction distribution deviates considerably from observed distribution for the  
 block of western the Netherlands, as reflected by the high value of mean CRPS. This is consistent  
 to the relatively low  $R^2$  observed for the same block. However, some blocks with relatively high  $R^2$   
 (in the north and south) have high CRPS. This indicates that the prediction mean is well-predicted  
 but not the prediction interval (too narrow).

Customised CV

414 There is a distinctive difference between model performance in areas close to traffic (i.e. *tr-hp*  
415 and *tr-lmp*) and far away from traffic (i.e. *far*). The INLA model outperformed other non-spatial  
416 methods in both *tr-hp* and *tr-lmp*, especially for the latter while the XGB model outperformed the  
417 INLA model (and all the other models) in *far*. This indicates the importance of modelling spatial  
418 dependency in areas close to traffic and possibly non-linear relationships far-away from roads. All the  
419 ensemble tree-based methods obtained much worse results compared to linear regression methods in  
420 *tr-lmp*. A linear regression model typically outperforms ensemble tree-based methods when there are  
421 relatively few observations for a flexible relationship to be justified. As the number of observations  
422 that are close to traffic and far away from traffic is balanced, the results indicate that the population  
423 density alters relationships between NO<sub>2</sub> and road density (i.e. the relationships between NO<sub>2</sub> and  
424 road density is different with different population density) in areas close to traffic.

## 425 5.2 Prediction interval

426 The 90% prediction intervals for INLA, INLA-G, DF, QRF and QRFLA are shown in figs. 4 to 6.  
427 The RF-based methods, namely DF, QRF and QRFLA reach the coverage probability higher than  
428 0.9, but the DF predicts a more realistic prediction quantile, notably, it covers four observations that  
429 are not covered by the same prediction quantiles predicted by the QRF. The INLA 90% prediction  
430 interval is too narrow. The coverage probability is 0.41 for INLA and 0.36 for INLA-G. The predicted  
431 90th quantile of the INLA-G turned to better capture extreme high values but the model also turned  
432 to miss more at lower values. The QRFLA predicted a slightly narrower prediction interval compared  
433 to QRF. This indicates that Lasso reduced the variance of a QRF model by aggregating trees.

Table 4: Results with customised CV. tr-hp: close to traffic and high population, tr-lmp: close to traffic and middle and low population, far: far away from traffic. RRMSE (relative RMSE), RMAE (relative MAE), RIQR (relative IQR).

	RMSE	RRMSE	IQR	RIQR	MAE	RMAE	$R^2$
LA_tr-hp	12.4	0.3	17.3	0.4	10.2	0.3	0.11
RF_tr-hp	11.9	0.3	17.8	0.5	9.8	0.3	0.18
XGB_tr-hp	11.6	0.3	15.3	0.4	9.3	0.2	0.21
INLA_tr-hp	11.3	0.3	16.6	0.4	9.5	0.3	0.26
LA_tr-lmp	7.5	0.3	10.4	0.5	6.1	0.3	0.21
RF_tr-lmp	8.2	0.4	10.9	0.5	6.4	0.3	0.05
XGB_tr-lmp	8.2	0.4	10.5	0.5	6.4	0.3	0.04
INLA_tr-lmp	6.7	0.3	8.7	0.4	5.3	0.2	0.36
LA_far	5.0	0.4	4.9	0.4	4.2	0.3	0.47
RF_far	4.9	0.3	4.0	0.3	3.6	0.3	0.47
XGB_far	3.4	0.2	3.6	0.3	2.5	0.2	0.74
INLA_far	4.0	0.3	4.3	0.3	3.2	0.2	0.65

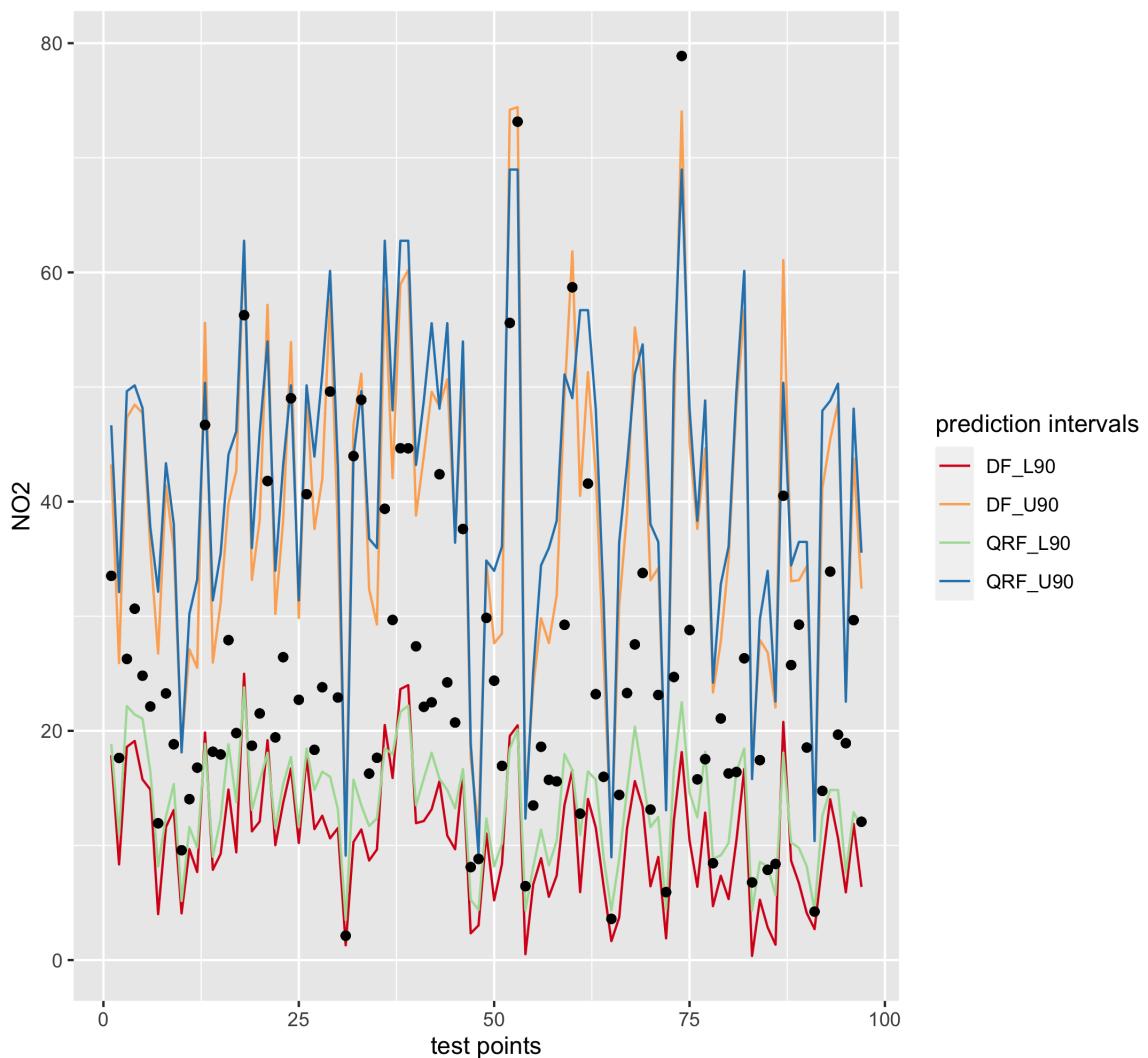


Figure 4: The 90% prediction interval predicted by DF and QRF. The black dots indicate observations in the test dataset.

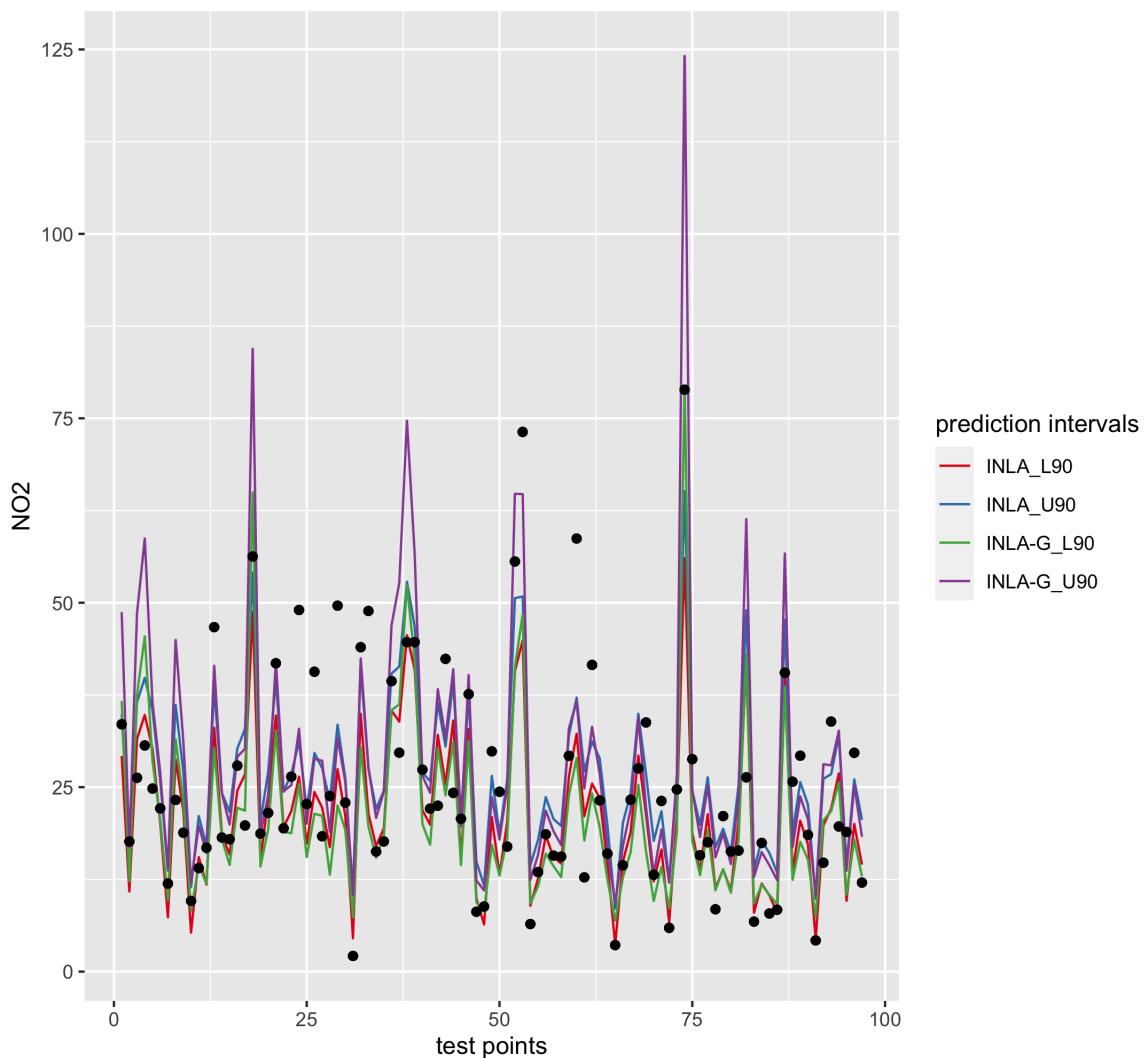


Figure 5: The 90% prediction interval predicted by INLA and INLA-G. The black dots indicate observations in the test dataset.

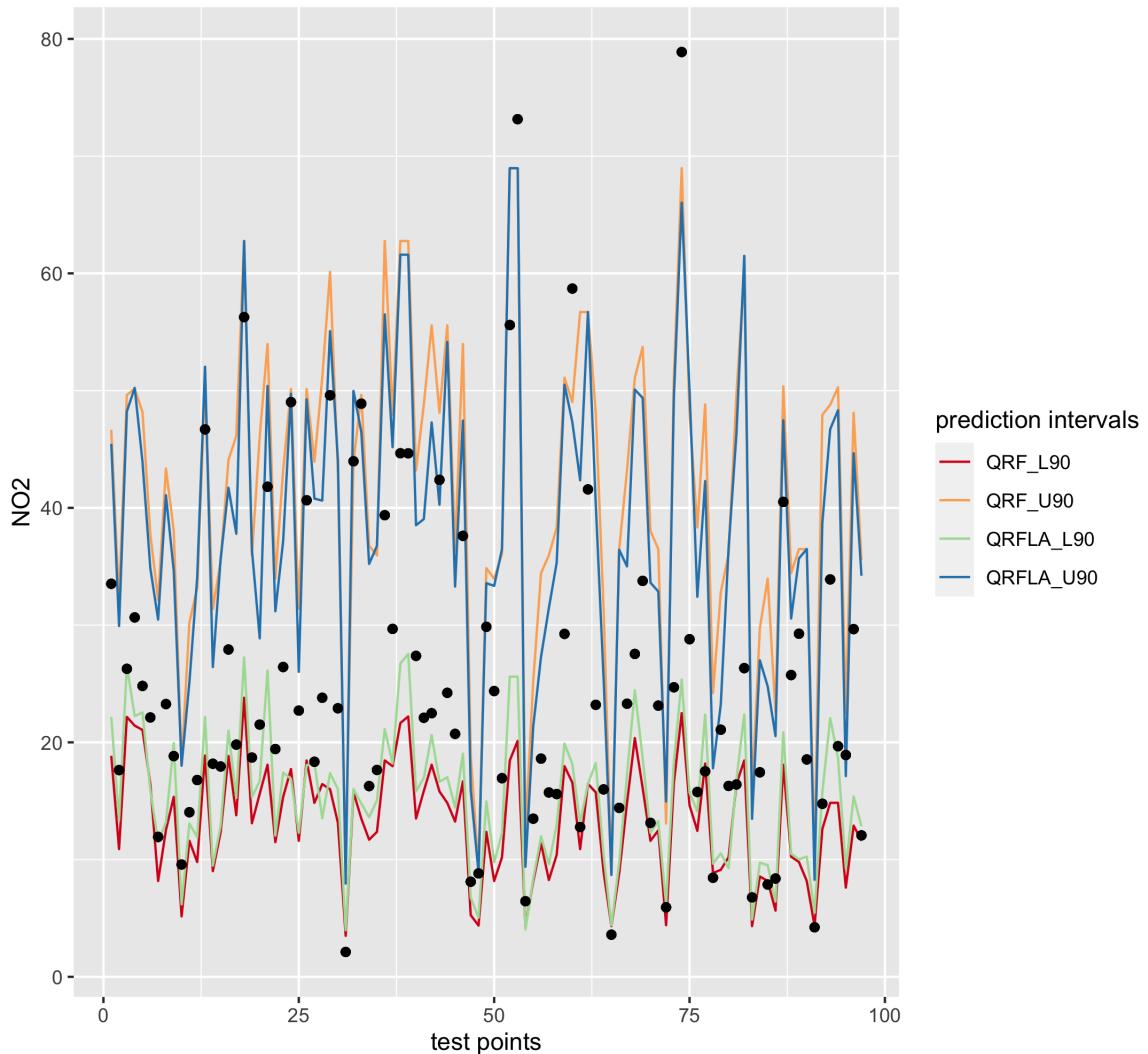


Figure 6: The 90% prediction interval predicted by QRF and QRFLA. The black dots indicate observations in the test dataset.

<sup>434</sup> **5.3 Model Interpretation**

<sup>435</sup> SHAP values are calculated for RF and XGB methods using all the data. The variables are ranked by  
<sup>436</sup> their variable importance, which is calculated as the sum of SHAP magnitudes over all the samples. It  
<sup>437</sup> can be observed from fig. 7 that the variable rankings and the pattern of variable impacts on model  
<sup>438</sup> output are similar. Both methods ranked road\_class\_2\_100 at the top. The variable importance  
<sup>439</sup> calculated by the SHAP indicates a pattern that matches well with our expectation in the emission  
<sup>440</sup> sources (e.g. high pollution close to primary roads). To illustrate, we observe a positive trend of  
<sup>441</sup> SHAP values along with road\_class\_2\_100 values, this matches with the explanation that areas with  
<sup>442</sup> higher primary road density generally experience higher NO<sub>2</sub> concentrations.

<sup>443</sup> To analyse the effect of each covariate in the INLA model, we firstly normalised all the covariates  
<sup>444</sup> (by subtracting the mean and dividing the centred columns by their standard deviations) and used  
<sup>445</sup> all the data to fit the INLA model. road\_class\_2\_100 has the highest effect (mean = 4.37), follows by  
<sup>446</sup> the population\_3000 (3.08), these are consistent to the XGB variable importance (fig. 7b ). Then,  
<sup>447</sup> the road\_class\_3\_300 (3.00) has a notably higher effect (besides the top 2) than other covariates,  
<sup>448</sup> which has coefficients from 0.72 to 1.88. This differs from the XGB and RF variable importance  
<sup>449</sup> which ranked the population\_1000 higher above, while in the INLA model the population\_1000 has  
<sup>450</sup> the lowest effect (0.72). This may be because of the high correlation between population\_1000 and  
<sup>451</sup> population\_3000, as SHAP is a permutation test, it ignores the dependency between covariates.  
<sup>452</sup> In general, both geostatistical and ML methods estimated covariate effects match their physical  
<sup>453</sup> explanations. The statistics (mean, standard deviation, mode) and predicted quantiles of each  
<sup>454</sup> coefficient are shown in the supplementary material figure 3.

<sup>455</sup> The differences between the predicted NO<sub>2</sub> and the mean of the spatial random field fig. 8  
<sup>456</sup> indicates the effects of covariates. The highest values of the mean of the spatial random field are  
<sup>457</sup> shown close to the Stuttgart region. Relatively high values can be observed in northern, southern

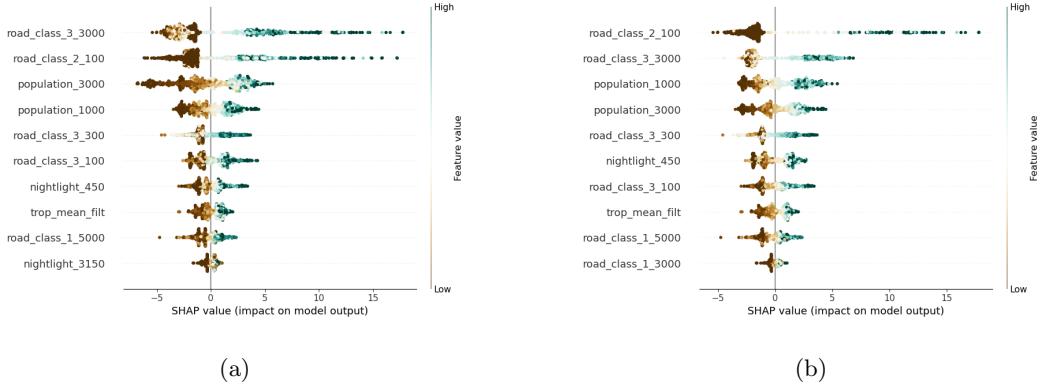


Figure 7: Variable impact calculated by SHAP (SHapley Additive exPlanations), a) the RF model, b) The XGB model. The horizontal location shows whether the effect of that value is associated with a higher or lower prediction. The covariate ranking is based on the sum of SHAP magnitudes over all the samples.

458 and western Germany. Compared to fig. 9, the areas close to the Stuttgart (Germany) region where  
 459 the mean values of the spatial random field are high corresponds to the high magnitudes of NO<sub>2</sub>  
 460 concentrations. Also, the differences between the observations and predictions are relatively large  
 461 in magnitudes in this region. To facilitate visualisation, we also calculated the differences between  
 462 INLA model predictions and the observations (supplementary material, figure 2).

## 463 6 Discussion

464 In this study, we compared geostatistical methods with ML methods for spatial NO<sub>2</sub> prediction in  
 465 Germany and the Netherlands. The comparison consists of the predicted mean, prediction inter-  
 466 vals, and model interpretation. Spatial and non-spatial CV strategies are used to reveal prediction  
 467 accuracy in different aspects. We also implemented the Lasso post-processed RF and geostatistical  
 468 stacked learning for NO<sub>2</sub> mapping (which to our knowledge have not been applied in air pollution

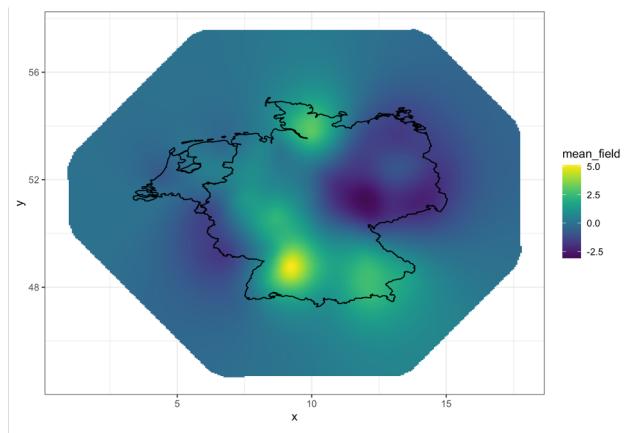


Figure 8: Mean of the spatial random field fitted by the INLA model.

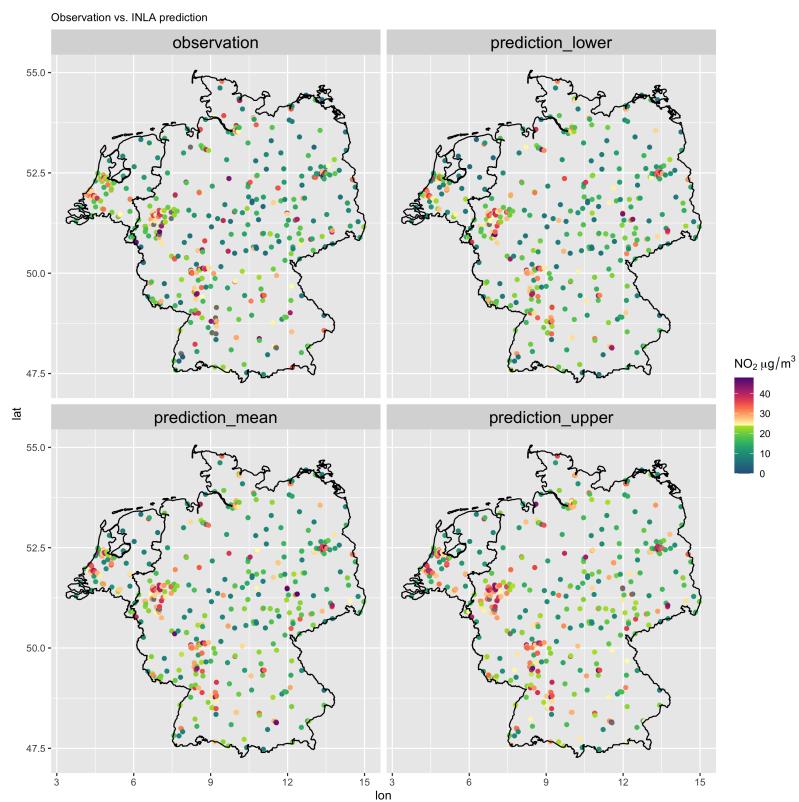


Figure 9: INLA predicted NO<sub>2</sub> at the ground stations with mean (prediction\_mean), high (prediction\_high, 0.975) and low (prediction\_low, 0.925) quantiles and the observed NO<sub>2</sub> (observation).

<sup>469</sup> mapping before) and these two methods considerably improve from the original RF and stacked  
<sup>470</sup> learning methods, respectively.

<sup>471</sup> Several venues were attempted to further improve the geostatistical model fitted with INLA.  
<sup>472</sup> Firstly, as we observed in general worse results at the geographical boundaries (figs. 2 and 3), we  
<sup>473</sup> inspected if different meshes with edge-effects fully accounted (e.g. the mesh is sufficiently large for  
<sup>474</sup> observations at the edge) could improve the prediction accuracy. It turned out that the same perfor-  
<sup>475</sup> mance is obtained. Secondly, we suspected that the deviation from assumed distribution (Gaussian)  
<sup>476</sup> is the cause of narrow prediction intervals of the INLA model. However, assuming a Gamma distri-  
<sup>477</sup> bution likelihood did not improve the model performance (in terms of the accuracy matrix, CRPS  
<sup>478</sup> and coverage probability). We also experienced the square transformation of the observations and  
<sup>479</sup> the use of the log-normal likelihood but that also decreases the model performance. Thirdly, we  
<sup>480</sup> additionally added two factor variables, namely "country code" (country code, "DE" for Germany  
<sup>481</sup> and "NL" for the Netherlands) and "urban types" (rural, urban, city centre according to (Dijkstra  
<sup>482</sup> and Poelman, 2014)). However, that also does not increase the model performance. In future works,  
<sup>483</sup> using a different spatial model (e.g. by specifying different hyperparameters), using the country and  
<sup>484</sup> urban types as mixed-effects, and modelling spatial varying coefficients may improve the modelling  
<sup>485</sup> results. Major improvement may also be achieved by integrating mobile sensing measurements and  
<sup>486</sup> other geospatial predictors (e.g. traffic count, urban morphological matrix) (Moraga et al., 2017).

<sup>487</sup> We implemented an INLA model without modelling the spatial random effect (called non-spatial  
<sup>488</sup> INLA) to deepen our understanding of the effect of modelling the spatial process in our INLA model.  
<sup>489</sup> The non-spatial INLA model obtained lower DIC (Information Criterion) 3286.66 vs. 3251.97 (with  
<sup>490</sup> spatial effects) and WAIC (Watanabe-Akaike information criterion) 3291.75 vs. 3253.93 (with spa-  
<sup>491</sup> tial effects). These suggest the advantage of modelling the spatial effects. We normalised covari-  
<sup>492</sup> ates before inputting into the spatial and non-spatial INLA models and compared the differences

493 between the fixed-effects obtained by the original and non-spatial INLA model (supplementary ma-  
494 terial figure 3-4) and found the most notable change is on the increased effect on the covariate  
495 population\_1000 for the non-spatial INLA model. This can be explained by that part of the effects  
496 of population\_1000 is modelled in the spatial random field. The second most notable change is on  
497 the decreased effect of nightlight\_450 for the non-spatial INLA model. After the spatial process is  
498 modelled, the nightlight\_450 has a higher contribution to the model. Together with the decreased  
499 effects of road\_class\_2\_100 and road\_class\_3\_300 for the non-spatial INLA model, these may indicate  
500 that the spatial model could better account for traffic-related variables (i.e. road and nightlight in  
501 smaller buffers).

502 Model performance differs between the three road and population situations. The "far" situation  
503 obtained the best modelling accuracy while the "tr-hp" the worst. This is likely due to the fact that  
504 the urban NO<sub>2</sub> process is more complex due to urban forms and traffic conditions. This may also  
505 indicate that more detailed traffic counts and meteorological data are needed for modelling the NO<sub>2</sub>  
506 emission sources.

507 Different from non-parametric models such as ensemble trees, a parametric geostatistical model  
508 fitted with INLA as the one developed in our study requires feature selection and the assumption  
509 of the distribution of the response. Several studies used the whole dataset for variable selection and  
510 then use selected variables for CV (Lu et al., 2020b; Larkin et al., 2017). This may however lead to  
511 an information leak as the validation data is also used in CV. To avoid this problem, one can include  
512 the variable selection process in each CV (i.e. use the same training data for variable selection and  
513 test). However, variable selection in each run added in additional error and uncertainty, therefore,  
514 a determined set of covariates may be preferred. We obtain a fixed set of selected variables while  
515 reducing information leakage to a negligible level by choosing only the variables that are selected  
516 90% -100% times of all the bootstraps of Lasso.

517 Using the geostatistical method to stack learners obtained higher prediction accuracy in terms  
518 of the mean prediction compared to the non-spatial stacking. This suggests the complex response-  
519 covariate relationships modelled by the ML learners do not fully capture the spatial process. The  
520 geostatistical stacked models obtained the highest prediction accuracy and with high-performance  
521 computation, it is possible to apply them to a large-scale and at a high resolution. The limitation of  
522 such stacked methods is that they cannot be used to analyse the effects of covariates and therefore  
523 NO<sub>2</sub> emission sources. But these models could be a reference to the level of accuracy a statistical  
524 predictive model could reach with the data available and the characteristics of the base learners  
525 (here: if the base learners are global or local models).

## 526 7 Conclusion

527 We proposed a model comparison process to comprehensively compare between models considering  
528 not only the predicted mean but also prediction intervals and model interpretation. We also showed  
529 that the information provided by commonly single-used non-spatial CV may miss reflecting model  
530 behaviours. With the model comparison process, we compared the use of geostatistical and ML  
531 methods for the spatial prediction of NO<sub>2</sub> in Germany and the Netherlands and found noticeable  
532 differences in their limitations and strength. The geostatistical models are preferred especially for  
533 urban area prediction and provide the spatial process of observations and indicate the insufficient  
534 modelling of spatial random-effects of fixed-effects. But the uncertainty assessment of geostatistical  
535 methods, which is commonly known as strength, fails to provide a prediction interval that meets  
536 the expectation. The QRF and DF obtained satisfying prediction intervals, with the DF slightly  
537 more capable of predicting the extremes. Using Lasso to aggregate trees in random forest increase  
538 model performance and reduce model variance. Using the geostatistical method to stack learners  
539 obtained the highest accuracy in terms of the mean prediction. Despite the NO<sub>2</sub> observations follow

<sup>540</sup> closer to a Gamma distribution than a Gaussian, the use of a Gamma likelihood in the geostatistical  
<sup>541</sup> model and Gamma objective in the XGBoost obtained much worse results than using a Gaussian  
<sup>542</sup> likelihood or objective. By comparing with the non-spatial stacking, geostatistical stacking suggests  
<sup>543</sup> the necessity of modelling the spatial process.

544 **References**

- 545 C. Alakus, D. Larocque, and A. Labbe. Rfpredinterval: An r package for prediction intervals with  
546 random forests and boosted forests. *arXiv preprint arXiv:2106.08217*, 2021.
- 547 L. Anselin et al. Spatial econometrics. *A companion to theoretical econometrics*, 310330, 2001.
- 548 A. Beloconi and P. Vounatsou. Bayesian geostatistical modelling of high-resolution no<sub>2</sub> exposure  
549 in europe combining data from monitors, satellites and chemical transport models. *Environment*  
550 *International*, 138:105578, 2020. ISSN 0160-4120. doi: <https://doi.org/10.1016/j.envint.2020.105578>. URL <https://www.sciencedirect.com/science/article/pii/S0160412019324109>.
- 552 S. Bertazzon, M. Johnson, K. Eccles, and G. G. Kaplan. Accounting for spatial effects in land use  
553 regression for urban air pollution modeling. *Spatial and Spatio-temporal Epidemiology*, 14-15:9 –  
554 21, 2015. ISSN 1877-5845.
- 555 S. Bhatt, E. Cameron, S. R. Flaxman, D. J. Weiss, D. L. Smith, and P. W. Gething. Improved  
556 prediction accuracy for disease risk mapping using gaussian process stacked generalization. *Journal*  
557 *of the Royal Society Interface*, 14(134):20170520, 2017.
- 558 M. Blangiardo and M. Cameletti. *Spatial and spatio-temporal Bayesian models with R-INLA*. John  
559 Wiley & Sons, 2015.
- 560 L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 561 D. J. Briggs, C. de Hoogh, J. Gulliver, J. Wills, P. Elliott, S. Kingham, and K. Smallbone. A  
562 regression-based method for mapping traffic-related air pollution: application and testing in four  
563 contrasting urban environments. *Science of the Total Environment*, 253(1-3):151–167, 2000.
- 564 J. Chen, K. de Hoogh, J. Gulliver, B. Hoffmann, O. Hertel, M. Ketzel, M. Bauwelinck, A. van  
565 Donkelaar, U. A. Hvidtfeldt, K. Katsouyanni, et al. A comparison of linear regression, regular-

566 ization, and machine learning algorithms to develop Europe-wide spatial models of fine particles  
567 and nitrogen dioxide. *Environment international*, 130:104934, 2019a.

568 T. Chen and C. Guestrin. xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm*  
569 *sigkdd international conference on knowledge discovery and data mining*, pages 785–794. ACM,  
570 2016.

571 T. Chen, T. He, M. Benesty, V. Khotilovich, Y. Tang, H. Cho, K. Chen, R. Mitchell, I. Cano,  
572 T. Zhou, M. Li, J. Xie, M. Lin, Y. Geng, and Y. Li. *xgboost: Extreme Gradient Boosting*, 2019b.  
573 URL <https://CRAN.R-project.org/package=xgboost>. R package version 0.82.1.

574 M. Chiusolo, E. Cadum, M. Stafoggia, C. Galassi, G. Berti, A. Faustini, L. Bisanti, M. A. Vigotti,  
575 M. P. Dessì, A. Cerniglio, et al. Short-term effects of nitrogen dioxide on mortality and sus-  
576 ceptibility factors in 10 italian cities: the epiair study. *Environmental health perspectives*, 119(9):  
577 1233–1238, 2011.

578 Copernicus. Sentinel-5p nrti no2: Near real-time nitrogen dioxide. [https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS\\_S5P\\_NRTI\\_L3\\_NO2#bands](https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S5P_NRTI_L3_NO2#bands), 2021. last ac-  
579 ccessed: Aug 3, 2021.

581 D. P. Dee, S. M. Uppala, A. Simmons, P. Berrisford, P. Poli, S. Kobayashi, U. Andrae, M. Balmaseda,  
582 G. Balsamo, d. P. Bauer, et al. The era-interim reanalysis: Configuration and performance of the  
583 data assimilation system. *Quarterly Journal of the royal meteorological society*, 137(656):553–597,  
584 2011.

585 P. J. Diggle, P. Moraga, B. Rowlingson, and B. M. Taylor. Spatial and spatio-temporal log-gaussian  
586 cox processes: extending the geostatistical paradigm. *Statistical Science*, 28(4):542–563, 2013.

587 L. Dijkstra and H. Poelman. *A harmonised definition of cities and rural areas: the new degree of*

- 588      *urbanisation*, 2014. URL [https://ec.europa.eu/regional\\_policy/sources/docgener/work/2014\\_01\\_new\\_urban.pdf](https://ec.europa.eu/regional_policy/sources/docgener/work/2014_01_new_urban.pdf). Last accessed: Aug 4, 2021.
- 590      T. Duan, A. Anand, D. Y. Ding, K. K. Thai, S. Basu, A. Ng, and A. Schuler. Ngboost: Natural  
591      gradient boosting for probabilistic prediction. In *International Conference on Machine Learning*,  
592      pages 2690–2700. PMLR, 2020.
- 593      Earthdata. *GES DISC*. URL "[https://disc.gsfc.nasa.gov/datasets/OMN02d\\_003/summary?keywords=OMI%202017%20No2](https://disc.gsfc.nasa.gov/datasets/OMN02d_003/summary?keywords=OMI%202017%20No2)". last assessed May 21, 2019.
- 594      EEA. *Explore air pollution data*, 2021. URL <https://www.eea.europa.eu/themes/air/explore-air-pollution-data>.
- 595      F. Fouedjio and J. Klump. Exploring prediction uncertainty of spatial data in geostatistical and  
596      machine learning approaches. *Environmental Earth Sciences*, 78(1):38, 2019.
- 597      J. H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*,  
598      pages 1189–1232, 2001.
- 599      E. Gonzalez-Estrada and J. A. Villasenor-Alva. *goft: Tests of Fit for some Probability Distributions*,  
600      2020. URL <https://CRAN.R-project.org/package=goft>. R package version 1.3.6.
- 601      T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning: data mining, inference,  
602      and prediction*. Springer Science & Business Media, 2009.
- 603      T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning: data mining, inference,  
604      and prediction, second edition*. Springer Science & Business Media, 2017.
- 605      G. Hoek, R. Beelen, K. De Hoogh, D. Vienneau, J. Gulliver, P. Fischer, and D. Briggs. A review  
606      of land-use regression models to assess spatial variation of outdoor air pollution. *Atmospheric  
607      environment*, 42(33):7561–7578, 2008.

- 610 G. James, D. Witten, T. Hastie, and R. Tibshirani. *An introduction to statistical learning*, volume  
611 112. Springer, 2013.
- 612 A. Jordan, F. Krüger, and S. Lerch. Evaluating probabilistic forecasts with scoringrules. *arXiv*  
613 *preprint arXiv:1709.04743*, 2017.
- 614 J. Kerckhoffs, G. Hoek, L. Portengen, B. Brunekreef, and R. C. Vermeulen. Performance of pre-  
615 diction algorithms for modeling outdoor air pollution spatial surfaces. *Environmental science &*  
616 *technology*, 53(3):1413–1421, 2019.
- 617 E. T. Krainski, V. Gómez-Rubio, H. Bakka, A. Lenzi, D. Castro-Camilo, D. Simpson, F. Lindgren,  
618 and H. Rue. *Advanced spatial modeling with stochastic partial differential equations using R and*  
619 *INLA*. CRC Press, 2018.
- 620 A. Larkin, J. A. Geddes, R. V. Martin, Q. Xiao, Y. Liu, J. D. Marshall, M. Brauer, and P. Hystad.  
621 Global land use regression model for nitrogen dioxide air pollution. *Environmental Science &*  
622 *Technology*, 51(12):6957–6964, 2017.
- 623 J. J. Li, A. Jutzeler, B. Faltings, S. Winter, and C. Rizos. Estimating urban ultrafine particle  
624 distributions with gaussian process models. *Research@ Locate14*, pages 145–153, 2014.
- 625 F. Lindgren, H. Rue, and J. Lindström. An explicit link between gaussian fields and gaussian  
626 markov random fields: the stochastic partial differential equation approach. *Journal of the Royal*  
627 *Statistical Society: Series B (Statistical Methodology)*, 73(4):423–498, 2011.
- 628 F. Lindgren, H. Rue, et al. Bayesian spatial modelling with r-inla. *Journal of Statistical Software*,  
629 63(19):1–25, 2015.
- 630 Y. Liu, G. Cao, and N. Zhao. Integrate machine learning and geostatistics for high-resolution

- 631 mapping of ground-level pm2. 5 concentrations. In *Spatiotemporal Analysis of Air Pollution and*  
632 *Its Application in Public Health*, pages 135–151. Elsevier, 2020.
- 633 M. Lu, O. Schmitz, K. de Hoogh, Q. Kai, and D. Karssenberg. Evaluation of different methods  
634 and data sources to optimise modelling of no2 at a global scale. *Environment international*, 142:  
635 105856, September 2020a. ISSN 1873-6750. doi: 10.1016/j.envint.2020.105856.
- 636 M. Lu, I. Soenario, M. Helbich, O. Schmitz, G. Hoek, M. van der Molen, and D. Karssenberg. Land  
637 use regression models revealing spatiotemporal co-variation in no2, no, and o3 in the netherlands.  
638 *Atmospheric Environment*, 223:117238, 2020b.
- 639 M. Lu, R. Dai, C. de Boer, O. Schmitz, I. Kooter, S. Cristescu, and D. Karssenberg. *Problems*  
640 *in Statistical Modelling of Air Pollution Basing Solely on Ground Monitor Stations and a Novel*  
641 *Mobile Sensing Instrument Solution*, 2021. submitted to Science of the Total Environment.
- 642 S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In  
643 I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and  
644 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Cur-  
645 ran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf>.
- 647 S. M. Lundberg, B. Nair, M. S. Vavilala, M. Horibe, M. J. Eisses, T. Adams, D. E. Liston, D. K.-W.  
648 Low, S.-F. Newman, J. Kim, et al. Explainable machine-learning predictions for the prevention  
649 of hypoxaemia during surgery. *Nature Biomedical Engineering*, 2(10):749, 2018.
- 650 K. Luo, R. Li, W. Li, Z. Wang, X. Ma, R. Zhang, X. Fang, Z. Wu, Y. Cao, and Q. Xu. Acute effects  
651 of nitrogen dioxide on cardiovascular mortality in beijing: an exploration of spatial heterogeneity  
652 and the district-specific predictors. *Scientific reports*, 6(1):1–13, 2016.

- 653 S. Martino and H. Rue. Implementing approximate bayesian inference using integrated nested laplace  
654 approximation: A manual for the inla program. *Department of Mathematical Sciences, NTNU,*  
655 *Norway*, 2009.
- 656 T. G. Martins, D. Simpson, F. Lindgren, and H. Rue. Bayesian computing with inla: new features.  
657 *Computational Statistics & Data Analysis*, 67:68–83, 2013.
- 658 N. Meinshausen. Quantile regression forests. *Journal of Machine Learning Research*, 7(Jun):983–999,  
659 2006.
- 660 P. Moraga. *Geospatial Health Data: Modeling and Visualization with R-INLA and Shiny*. Chapman  
661 & Hall/CRC, 2019.
- 662 P. Moraga, S. M. Cramb, K. L. Mengersen, and M. Pagano. A geostatistical model for combined  
663 analysis of point-level and area-level data using inla and spde. *Spatial Statistics*, 21:27–41, 2017.
- 664 NASA. *Shuttle Radar Topography Mission*. URL [https://www2.jpl.nasa.gov/srtm/  
665 dataprelimdescriptions.html](https://www2.jpl.nasa.gov/srtm/dataprelimdescriptions.html). last assessed Aug 15, 2021.
- 666 D. A. Nelson. European environment agency. *Colo. J. Int'l Envtl. L. & Pol'y*, 10:153, 1999.
- 667 NOAA. Dmsp and viirs data download. "<https://ngdc.noaa.gov/eog/download.html>", 2021.  
668 Last Accessed: 11.03.2021.
- 669 OpenStreetMap contributors. Planet dump 7 Jan 2019 retrieved from <https://planet.osm.org>, 2019.
- 670 X. Ren, Z. Mi, and P. G. Georgopoulos. Comparison of machine learning and land use regression  
671 for fine scale spatiotemporal estimation of ambient air pollution: Modeling ozone concentrations  
672 across the contiguous united states. *Environment International*, 142:105827, 2020. ISSN 0160-  
673 4120. doi: <https://doi.org/10.1016/j.envint.2020.105827>. URL <https://www.sciencedirect.com/science/article/pii/S0160412020317827>.

- 675 H. Rue and L. Held. *Gaussian Markov random fields: theory and applications*. CRC press, 2005.
- 676 H. Rue, S. Martino, and N. Chopin. Approximate bayesian inference for latent gaussian models by  
677 using integrated nested laplace approximations. *Journal of the royal statistical society: Series b*  
678 (*statistical methodology*), 71(2):319–392, 2009.
- 679 Y. Rybarczyk and R. Zalakeviciute. Machine learning approaches for outdoor air quality modelling:  
680 A systematic review. *Applied Sciences*, 8(12):2570, 2018.
- 681 L. Schlosser, T. Hothorn, R. Stauffer, A. Zeileis, et al. Distributional regression forests for prob-  
682 abilistic precipitation forecasting in complex terrain. *The Annals of Applied Statistics*, 13(3):  
683 1564–1589, 2019.
- 684 G. Shaddick, M. L. Thomas, H. Amini, D. Broday, A. Cohen, J. Frostad, A. Green, S. Gumy, Y. Liu,  
685 R. V. Martin, et al. Data integration for the assessment of population exposure to ambient air  
686 pollution for global burden of disease assessment. *Environmental science & technology*, 52(16):  
687 9069–9078, 2018.
- 688 D. M. Stasinopoulos, R. A. Rigby, et al. Generalized additive models for location scale and shape  
689 (gamlss) in r. *Journal of Statistical Software*, 23(7):1–46, 2007.
- 690 M. L. Stein. *Interpolation of spatial data: some theory for kriging*. Springer Science & Business  
691 Media, 2012.
- 692 J. A. Suykens and J. Vandewalle. Least squares support vector machine classifiers. *Neural processing*  
693 *letters*, 9(3):293–300, 1999.
- 694 J. Velthoen, C. Dombry, J.-J. Cai, and S. Engelke. Gradient boosting for extreme quantile regression.  
695 *arXiv preprint arXiv:2103.00808*, 2021.

- 696 A. M. Vicedo-Cabrera, A. Biggeri, L. Grisotto, F. Barbone, and D. Catelan. A bayesian kriging  
697 model for estimating residential exposure to air pollution of children living in a high-risk area in  
698 italy. *Geospatial health*, 8(1):87–95, 2013.
- 699 J. A. Villaseñor and E. González-Estrada. A variance ratio test of fit for gamma distributions.  
700 *Statistics & Probability Letters*, 96:281–286, 2015.
- 701 S. Wager, T. Hastie, and B. Efron. Confidence intervals for random forests: The jackknife and the  
702 infinitesimal jackknife. *The Journal of Machine Learning Research*, 15(1):1625–1651, 2014.
- 703 Q. Wang, H. Feng, H. Feng, Y. Yu, J. Li, and E. Ning. The impacts of road traffic on urban air  
704 quality in jinan based gwr and remote sensing. *Scientific Reports*, 11(1):1–9, 2021.
- 705 M. T. Young, M. J. Bechle, P. D. Sampson, A. A. Szpiro, J. D. Marshall, L. Sheppard, and J. D.  
706 Kaufman. Satellite-based no<sub>2</sub> and model validation in a national prediction model based on  
707 universal kriging and land-use regression. *Environmental science & technology*, 50(7):3686–3694,  
708 2016.
- 709 C. Yuan. Models and methods for computationally efficient analysis of large spatial and spatio-  
710 temporal data. 2011.
- 711 L. Zhai, S. Li, B. Zou, H. Sang, X. Fang, and S. Xu. An improved geographically weighted regression  
712 model for pm2. 5 concentration estimation in large areas. *Atmospheric Environment*, 181:145–154,  
713 2018.
- 714 Y. Zhan, Y. Luo, X. Deng, K. Zhang, M. Zhang, M. L. Grieneisen, and B. Di. Satellite-based  
715 estimates of daily NO<sub>2</sub> exposure in China using hybrid random forest and spatiotemporal kriging  
716 model. *Environmental science & technology*, 52(7):4180–4189, 2018.

<sup>717</sup> B. Zou, Q. Pu, M. Bilal, Q. Weng, L. Zhai, and J. E. Nichol. High-resolution satellite mapping  
<sup>718</sup> of fine particulates based on geographically weighted regression. *IEEE Geoscience and Remote*  
<sup>719</sup> *Sensing Letters*, 13(4):495–499, 2016.