

# Evaluation of Multi-variate Time Series Clustering for Imputation of Air Pollution Data

Wedad Alahamade<sup>1,3</sup>, Iain Lake<sup>2</sup>, Claire E. Reeves<sup>2</sup>, and Beatriz De La Iglesia<sup>1</sup>

<sup>1</sup>School of Computing Sciences, University of East Anglia, Norwich NR4 7TJ, UK

<sup>2</sup>School of Environmental Sciences, University of East Anglia, Norwich NR4 7TJ, UK

<sup>3</sup>School of Computing Sciences, Taibah University, Medina 42353, Saudi Arabia

**Correspondence:** Wedad Alahamade (W.Alahamade@uea.ac.uk)

**Abstract.** Air pollution is one of the world's leading risk factors for death, with 6.5 million deaths per year worldwide attributed to air pollution-related diseases. Understanding the behaviour of certain pollutants through air quality assessment can produce improvements in air quality management that will translate to health and economic benefits. However problems with missing data and uncertainty hinder that assessment.

We are motivated by the need to enhance the air pollution data available. We focus on the problem of missing air pollutant concentration data either because a limited set of pollutants is measured at a monitoring site or because an instrument is not operating, so a particular pollutant is not measured for a period of time.

In our previous work, we have proposed models which can impute a whole missing time series to enhance air quality monitoring. Some of these models are based on a Multivariate Time Series (MVTS) clustering method. Here, we apply our method to real data and show how different graphical and statistical model evaluation functions enable us to select the imputation model that produces the most plausible imputations. We then compare the Daily Air Quality Index (DAQI) values obtained after imputation with observed values incorporating missing data. Our results show that using an ensemble model that aggregates the spatial similarity obtained by the geographical correlation between monitoring stations and the fused temporal similarity between pollutants concentrations produced very good imputation results. Furthermore, the analysis enhances understanding of the different pollutant behaviours, and of the characteristics of different stations according to their environmental type.

## 1 Introduction

Time Series (TS) analysis has received much attention in recent decades due to its importance in many real-world applications such as earthquake prediction (Di Bello et al., 1996), weather forecasting (Carbajal-Hernández et al., 2012), air pollution forecasting (Du et al., 2020), or human activity recognition (Seto et al., 2015). Generally speaking, TS data can be described as a sequence of observations that a variable takes over time. When several variables are observed and recorded simultaneously, this becomes a Multivariate Time Series (MVTS).

The quality of the air in the UK is assessed based on five main pollutants. In this study we focus on the four main pollutants; particulate matter less than 2.5 m in diameter ( $PM_{2.5}$ ) or less than 10 m in diameter ( $PM_{10}$ ), ozone ( $O_3$ ), and nitrogen dioxide ( $NO_2$ ). These pollutants are measured hourly at various monitoring stations.

The main challenge with analysing these pollutant TS is that not all the stations report all the pollutants. Even if a station does, it may not measure a particular pollutant all the time due to instrument down-time. In our previous work (Alahamade et al., 2021), we applied an intermediate fusion approach to fuse the distance between stations using the similarity of the four pollutants. The similarity between pollutant TS was measured using Shape-Based Distance (SBD) between hourly pollutant concentrations (TS), as we found that SBD is better than other measures on our dataset (Alahamade et al., 2020). Then we used the k-means clustering algorithm to cluster the stations based on the fused distance, we called that MVTS clustering. Our initial clustering analysis showed that using the basic k-means with the fused distance gives very compact geographical clustering that enhances our understanding of the UK's air pollutants behaviours. Adding to that, using the fused distance to measure the similarity between the pollutants helped us solve some of the uncertainty problems associated with missing pollutant values as the MVTS clustering enables imputation even when no measurement is available for a given pollutant. This is because the multivariate nature of the clustering enabled a station to be allocated to a cluster based on the value of the other pollutants measured.

Based on the clustering results and station geographical location, we proposed three models to impute the whole time series for the missing pollutant at a given station. In this paper, we apply multiple model evaluation functions to assess which model gives best results and to demonstrate the validity of our models.

Our long-term goal is to reduce the uncertainty in air quality assessment by imputing all missing pollutants in the monitoring stations. This will allow us to calculate new air quality indices that may/may not agree with the previous indices, that is the observed indices that incorporate missing data. This in turn will help us to identify where more measurements can be beneficial.

We refer to our approach as time series imputation because we used the observed time series to impute missing time series (whole TS) in stations where one pollutant is not measured but other pollutants are. In this process, we are not filling the missing values within the time series (e.g. interpolating) but imputing a new TS. Also, we do not use predictive models hence we do not consider this a prediction task. However, it could be argued that our task is close to spatial interpolation (Lam, 1983) even though it is not completely based on spatial information, that is, we did not use any geographical information within the proposed MVTS time series clustering. Geographical information, however, is used in the Nearest Neighbour approaches, which are used in the ensemble proposed. Nevertheless, the main goal of the spatial interpolation is to fill in the gaps (points/locations with unknown measurements) using points with known values to cover a certain geographical area (Lam, 1983). Our goal is to impute unmeasured pollutants (whole TS) in several stations where they are not measured using the fused similarity between stations of other pollutants or using an ensemble of techniques including the MVTS clustering approach. We would argue that our imputation approach incorporates some uncertainty by using a combination of values (within the clustering process and within the ensemble) to produce the imputed value.

The paper's structure is as follows: Section 2 discusses some of the existing TS clustering methods and their application in the air quality field. Section 3 gives a brief introduction of the air quality assessment in the UK and its challenges. Section

4 discusses in detail all the methods we used to impute the missing pollutants and evaluate our proposed solutions. Finally, in Section 5, we analyse the results of our imputation models. Then, we conclude the work with some final remarks and  
60 indication for further developments in Section 6.

## 2 Related work

In this section, we briefly review some representative research in clustering techniques and its application in air pollution modelling. Data mining techniques have been widely applied to study the air pollution data; however, most of this research focuses only on a single pollutant (univariate TS), while clustering multivariate time series remains a challenging task (Liao,  
65 2005). Partitioning algorithms such as k-means and k-medoids are very common among works related to TS clustering and have been applied in many papers (e.g. Ignaccolo et al. (2008); Austin et al. (2013); Tuysuzoglu et al. (2019))

Austin et al. (2013) used the k-means algorithm to identify spatial patterns in air pollution data to cluster the USA cities based on the similarity of their  $PM_{2.5}$  composition profiles, then characterise these clusters based on chemical characteristics, emission profiles, geographic locations and population density. Ignaccolo et al. (2008) transformed the TS of pollutant daily  
70 observations into a functional form to smooth the TS, then classified the air quality monitoring network in Northern Italy using the Partitioning Around Medoids algorithm (PAM) to cluster three individual pollutants, namely  $NO_2$ ,  $PM_{10}$ , and  $O_3$ . Tuysuzoglu et al. (2019) applied different clustering algorithms such as k-means, Expectation Maximisation, and Canopy for each air pollutants in the dataset ( $NO$ ,  $NO_2$ ,  $SO_2$ ,  $PM_{10}$ , and  $O_3$ ), then aggregated the clustering results based on majority voting to identify one clustering solution for similar regions in terms of air quality.

75 On the other hand, there has been some research into similarity within MVTS. For example, Fontes and Budman (2017) proposed a MVTS clustering method based on extracted features from the univariate TS. In their work, Principal Component Analysis (PCA) is used to measure the similarity between MVTS, and fuzzy k-means is used to cluster these TS. This clustering approach was used for fault detection in a gas turbine. Zhou and Chan (2014) developed an algorithm for clustering MVTS by discovering each TS's temporal patterns. Their algorithm is based on k-means and aims to groups MVTS with similar temporal  
80 patterns together into the same cluster. D'Urso et al. (2018) proposed robust fuzzy clustering models for MVTS based on an exponential transformation of the dissimilarities. This algorithm was applied to real-world data on the concentrations of three pollutants ( $NO$ ,  $NO_2$ , and  $PM_{10}$ ) in the Metropolitan City of Rome for the problem of detecting pollution alarms.

In our previous work (Alahamade et al., 2020), we compared different TS distance measures and imputation techniques to impute the missing observations and missing pollutants (TS). We found that using Shape-Based Distance (SBD) gives better  
85 separated cluster than Dynamic Time Warping (DTW). Also, using MICE to impute the TS missing observations is better than using some single imputation methods such as Simple Moving Average (SMA). We used a univariate TS clustering using k-medoids (PAM) to cluster stations and imputed the missing pollutants using the cluster average. In this work, we use the k-means clustering algorithm and include a number of pollutants in the clustering, which make it MVTS clustering. This clustering algorithm was proposed in Alahamade et al. (2021) where more details can be found. Here we extend that work

90 by applying the imputation solution to real data and using extensive evaluation methods to demonstrate its effectiveness. This enables us to extend our understanding of pollutant behaviour.

### 3 Air Quality Assessment

We will study air pollution using the concentrations measured at the Automatic Urban and Rural Network (AURN) around the UK. The stations in the network are automatic and produce hourly pollutant concentrations. The data is collected and stored, 95 then made directly available via the Web (DEFRA , 2021). There are 167 stations with different environmental types: rural, urban, suburban background, roadside, and industrial.

The Daily Air Quality Index (DAQI) represents air pollution levels in the UK. This index is reported based on the highest individual DAQI derived for each of the five major air pollutants ( $O_3$ ,  $NO_2$ ,  $PM_{10}$ ,  $PM_{2.5}$ , and  $SO_2$ ) based on their concentrations. If concentration data for some of these pollutants is not available, the DAQI is based on those pollutants for which 100 data is available. The DAQI is used to provide an indication of the air quality, and some associated information that may be used by at-risk groups as well as the general population (DEFRA , 2021). The DAQI is numbered from 1 to 10, and divided into four bands; ‘low’ (1–3), ‘moderate’ (4–6), ‘high’ (7–9) and ‘very high’ (10). The air quality is negatively correlated with DAQI index, meaning that a higher DAQI index represents worse air quality.

### 4 Methods

105 The MVTS clustering algorithm and our proposed imputation models were implemented in R, Version (3.5.2) and are fully explained in previous work (Alahamade et al., 2021). To provide a more robust testing scenario, we separate the ‘model building’ stage from the imputation testing stage. We use an initial data period of three years (2015-2017) as a training set to build the clustering, and then impute on the next year (2018) of the TS to evaluate the goodness of fit.

#### 4.1 Imputation Models of Missing Pollutant TS

110 For evaluation purposes, we assume each pollutant from each station is missing entirely and impute it. For any given station,  $j$ , to impute the values of missing pollutant  $P_i^j$ , where  $i$  represents the different pollutants ( $1 \leq i \leq 4$ ), we use different models under two main similarity criteria: the similarity using clustering solutions and the similarity using geographical distance.

The k-means clustering algorithm is used to group the stations based on their temporal similarity, which is the similarity in time between the hourly pollutant concentrations using SBD as the temporal distance measure. This distance function is 115 implemented in ‘dtwclust’ Package in R (Sarda-Espinosa et al., 2017). The geographical distance is used to find the spatial similarity between station locations. Adding to that, we use an ensemble model which calculates the median of all the previous imputation models; this model aggregates the temporal and the spatial imputation using both the time series clustering and the geographical location similarity. Then, we evaluate these models to select the one that gives the highest similarity to the real values which are known. We explain these models in detail in the following sections:

#### 120 4.1.1 Imputation models using clustering results:

Once a clustering of our stations is obtained, we can use the clustering solution to impute missing TS (pollutants). If station  $j$  belongs to cluster  $C_x$ , ( $1 \leq x \leq k$ , where  $k$  is the number of clusters) given the measured pollutants over time, then, to impute pollutant  $P_i$  based on the clustering results, we use three models:

- 125 1. We impute the average of pollutant  $P_i$  in cluster  $C_x$ , which is the hourly average of pollutant  $P_i$  in all the stations that fall in this cluster. We call this method Cluster Average (CA).
2. We impute the average of pollutant  $P_i$  in cluster  $C_x$ , but using only stations with the same environment type to station  $j$  within the cluster, such as ‘Background Rural’, ‘Background Urban’, ‘Traffic’, or ‘Industrial’. We call this method (CA+ENV). This is in recognition that the type of station may be important and result more similar pollutant concentrations.
- 130 3. We impute the average of pollutant  $P_i$  in cluster  $C_x$ , for stations that belong to the same region. As defined by Defra (DEFRA, 2021) there are 16 regions in the UK for air quality assessment, such as Eastern, North Wales, East Midlands, and the other UK regions; this method is called (CA+REG).

#### 4.1.2 Imputation models by similarity using geographical distance:

135 First, we measure the geographic distance using Harvison metric, which calculates geographic distance on earth based on longitude and latitude. We calculate the distance between station  $j$  and all other stations that measure pollutant  $P_i$ . Then to impute pollutant  $P_i$  for station  $j$  we use:

1. The nearest neighbour (1NN) using the Harvison based distance to station  $j$ ; this method is called (1NN).
2. The average of the two nearest neighbours (2NN) to station  $j$ ; this method is called (2NN).

#### 4.1.3 Imputation model by ensemble:

140 In this approach, for a given station  $j$ , to impute pollutant  $P_i$ , we use the median value of all the imputed values from the previous models. Those are cluster average (CA), cluster average considering the station type (CA+ENV), cluster average considering the region (CA+REG), first nearest neighbour (1NN), and the average of the two nearest neighbours (2NN). This method is called (Median). This imputation approach may be computationally the most expensive as it needs for all others to be computed, but ensembles have the potential to provide very powerful solutions by combining predictions.

#### 145 4.2 Imputation model evaluation

We evaluate how plausible the imputation is using different models by comparing truth values to imputed values. The models evaluation are based on the test dataset, which is the 2018 data. As earlier mentioned we do this by taking each existing TS for which we have values, one at a time, and consider them missing. We impute the whole TS by various models and

compare that to the ground truth. We are evaluating our models against the real concentrations which contain missing values,  
150 hence, we ignore all the missing values in this evaluation. For each model, we can average the different imputation models' behaviour from all the stations to establish the one that provides imputed values closest to the real values. Hence, for our experimental set up we take each existing TS for a given pollutant and station,  $P_i^j$  in turn, and impute it by the various models to obtain an imputed TS,  $PI_i^j$ . We compare the real values to the imputed values using different statistical and graphical model evaluation functions. The statistics function include Fraction of predictions within the factor of two (FAC2), Mean Bias (MB),  
155 Normalised Mean Bias (NMB), Root Mean Squared Error (RMSE), Coefficient of correlation (R), and Index of Agreement (IOA). These measures are used to evaluate the temporal variation of air pollutants between imputed/modelled and observed concentrations. The graphical functions include Conditional Quantile plot, Time Variation plot, and Taylor's Diagram. These are functions within the Openair Package, a freely available air quality data analysis tool in R (Carslaw and Ropkins, 2012), that present comparisons between the modelled and measured air pollutant concentrations and their statistics graphically. We  
160 use R packages 'openair' (Carslaw and Ropkins, 2012), and 'tidyverse' (Wickham and Wickham, 2017) for the evaluation.

Model evaluation functions are beneficial when more than one model is involved in the comparison, and help us in understanding why a model does not perform well. The model that gives the lowest error on average, the highest correlation and the highest degree of agreement between imputed and observed concentrations for all stations (i.e. imputed TS) is initially considered the best model. However, extensive evaluation with various graphical functions enable us to much better assess the  
165 model quality and how it reflects uncertainty. Note that the best model may change from one pollutant to another and may be affected by other factors such as station type (e.g. urban background, rural and roadside) or pollutant lifetime and spread.

### 4.3 DAQI calculation

In the UK, DAQI forecasts are issued on a national scale; they are produced by the Met Office in the morning for the current day as well as for the next four days. The forecast is improved by incorporating the recent observations of air quality recorded at  
170 the AURN stations. The overall air pollution index for a site or region is determined by the highest DAQI of the five pollutants. The regional DAQI is the highest index among all the stations at that region.

For our evaluation, we calculated the daily DAQI value using the observed data for each station. This is because the DAQI value is not saved as part of the historical data available so we need to calculate it from the downloaded data. Defra has published a guide for the implementation of DAQI (DEFRA , 2013), which explains how the value is calculated and we follow  
175 that guidance. To calculate DAQI, each air pollutant is calculated as follows:

- Ozone: the  $O_3$  is measured hourly. To determine the DAQI we need to calculate the daily maximum 8-hourly running mean concentration. First, for each hour we calculate the running 8-hourly mean from the previous hours. Then we find the maximum value of these 8-hourly running means. For this calculation 75% of the data must be captured to calculate the 8-hourly mean.
- 180 – Nitrogen dioxide: the  $NO_2$  is measured based on hourly mean. We calculate the daily  $NO_2$  contribution to the DAQI by taking the maximum observation in 24 hours every day from 0:00 to 23:00.

- Particles  $PM_{10}$   $PM_{2.5}$ : are measured hourly. The DAQI is based on the 24 hours mean, which we calculate by taking the mean value from the hourly observations. For these pollutants 75% of the daily observations must be captured to calculate the mean, otherwise, the pollutant is considered as missing that day.
- 185 – We define the daily index for each pollutant separately. Then, for a station, we take the highest air pollutant index to be the value of the DAQI at that station.

We called the DAQI that is calculated based on observation ‘observed DAQI’, and the DAQI that is calculated based on imputation ‘imputed DAQI’. We use the observed DAQI as a performance tool to evaluate our imputation model on its ability to reproduce the daily air quality index. Note that although we produce only one imputation and not multiple imputations at  
190 this stage, we believe they reflect the underlying uncertainty because they are based on a number of aggregated methods

## 5 Results

In this section, we first analyse the proposed pollutant imputation models using some statistical and graphical air pollution modeling evaluation functions. Then, we evaluate the imputation model performance based on the comparison between the observed and imputed DAQI.

### 195 5.1 Air pollution imputation modeling evaluation

We first evaluate imputation models based on the statistical and then on the graphical analysis.

#### 5.1.1 Model evaluation based on statistical analysis

Table. 1, shows the statistical analysis results. In this table  $N$  is the number of stations that measure each pollutant. The table also shows the Fraction of predictions within the factor of two (FAC2), Mean Bias (MB), Normalised Mean Bias (NMB), Root  
200 Mean Squared Error (RMSE), Coefficient of correlation (R), and Index of Agreement (IOA).

In general, model 6 (Median), which is the model that uses the ensemble technique of other models, gives the lowest error average (RMSE), the highest Pearson correlation coefficient (R), and the highest agreement between imputed and observed concentrations (IOA) for  $O_3$ ,  $PM_{2.5}$ , and  $PM_{10}$ . However,  $NO_2$  shows different behaviour with model 2 (CA+ENV) achieving slightly higher performance with an increase of the correlation coefficient (by 0.049) and decrease of error average (by 0.826)  
205 compared to model 6 (Median). The Model Bias (MB) for model 2 is 50% higher than that of model 6.  $NO_2$  shows local patterns, as it is concentrated where it is emitted in urban areas and near to the roadside. Adding to that  $NO_2$  is shorter lived than other pollutants and shows greater spatial variability, with concentrations being strongly influenced by the environment type (e.g. roadside, urban background, rural). This changes the  $NO_2$  concentrations from one location to another based on the environmental type (CenterForCities, 2020).

210 All the selected models performed well, with 71-89 % of their imputations falling within a factor of two of the observed concentrations as shown in the FAC2 values in Table. 1. According to Dick Derwent and Murrells (2010), an air quality model

minimum requirement is that the FAC2 value is higher than 0.50 and NMB values should be in the range between -0.2 and +0.2. Both are met by our models. NMB measures if the model under or over predict, as it estimates the difference between the mean observed and imputed concentrations. Negative NMB means that the model under-predict and vice versa. All the models  
215 have very small biases.

### 5.1.2 Model evaluation based on Taylor's diagram analysis

We use Taylor's diagram to analyse three main statistics: correlation coefficient R, the standard deviation ( $\sigma$ ) and the root-mean-square error (centred). These statistics can be plotted on one (2D) graph which can be represented through the Law of Cosines (Taylor, 2001).

220 The standard deviation represents the variability between modelled and observed concentrations. The observed variability is plotted on the x-axis. The magnitude of the variability is measured as the radial distance from the plot's origin. The black dashed line shows this for the observed value. The grey lines are isopleths for the correlation coefficient (R) as indicated by the arc shaped axis; the correlation increases along the arc towards the x-axis. The centred root-mean square error (RMS) is represented by the concentric brown dashed lines. The furthest the points/models are from the observed value the worst  
225 performance they have (Carslaw and Ropkins, 2012). Fig. 1 shows Taylor Diagram plots for all models with all pollutants.

In almost all cases the models exhibit less variability than the observed, indicated by the points being closer to the origin than the black dashed line. In general, Model 4 (INN) followed by Model 5 (2NN) show variability that is most similar to the observations as indicated by their relative closeness to the black dashed line. However, these models tend to have the lowest correlation coefficients, indicated by the grey lines, and the greatest RMSE, indicated by the brown dashed lines. Models 4 and  
230 5 use the concentrations from a single site (i.e. the nearest stations) in the imputation, where as the other models use a cluster average (CA, CA+REG, CA+ENV) or a model ensemble average (Median), so it is reasonable for model 4 and 5 to have fairly similar variability to the observed concentrations. All the other models display less variability than the observed concentrations (as indicated by their points being further from the black dashed line) which may be consistent with their derivation methods which may smooth out some of the variability.

235 Model 6 (Median), regardless of its ability to capture variability, is confirmed as having the highest correlation coefficient, and the lowest centred root means squared with all the pollutants except  $\text{NO}_2$ , for which it is the second best behind Model 2 (CA+ENV).

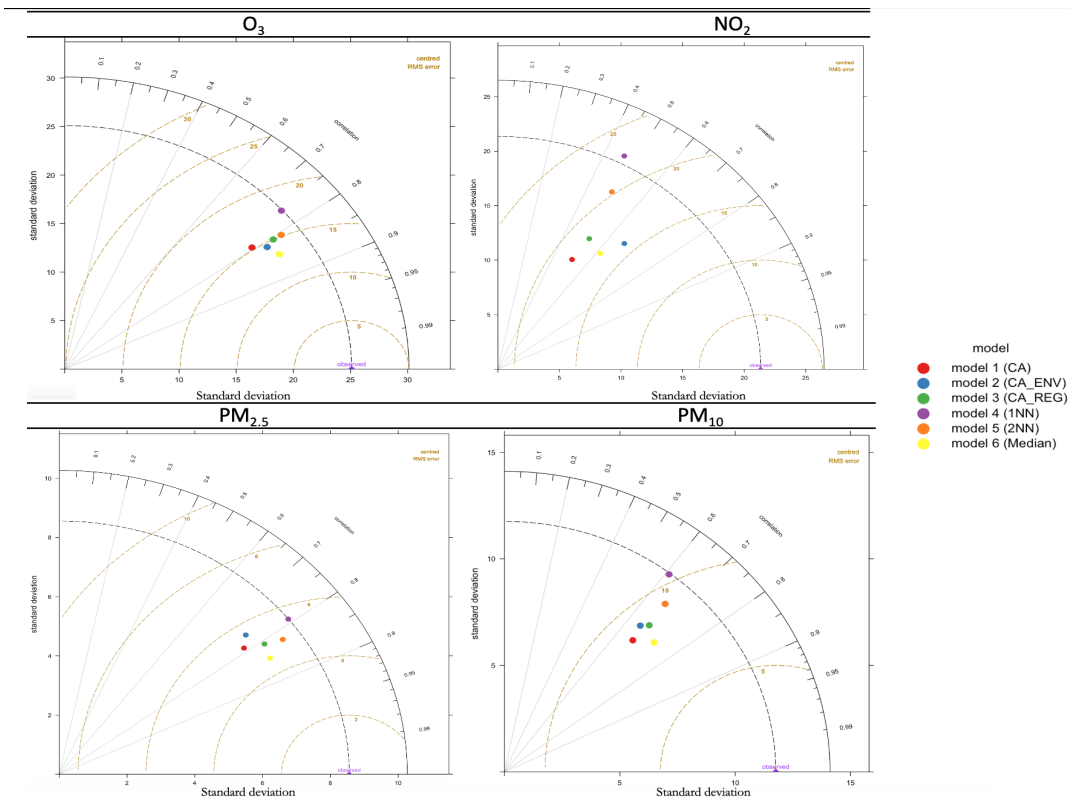
### 5.1.3 Model evaluation based on conditional quantile analysis

We analyse the spread of the modelled and observed pollutant concentrations using conditional quantile plots. Fig. 2 and 3 show  
240 the conditional quantile plots for the six imputation models (panels A to F). This visualisation splits the concentrations into bins according to values of the modelled concentrations. The median line of these values and the 25/75th and 10/90th quantile values are plotted together with a blue line showing a "perfect" model. Also shown are histograms of modelled concentrations (shaded grey bars) and histograms of observed concentrations (blue outline bars).



**Table 1.** Performance of the hourly pollutant concentrations imputation models based on statistical measures. Best values in bold for FAC2, RMSE, R and IOA

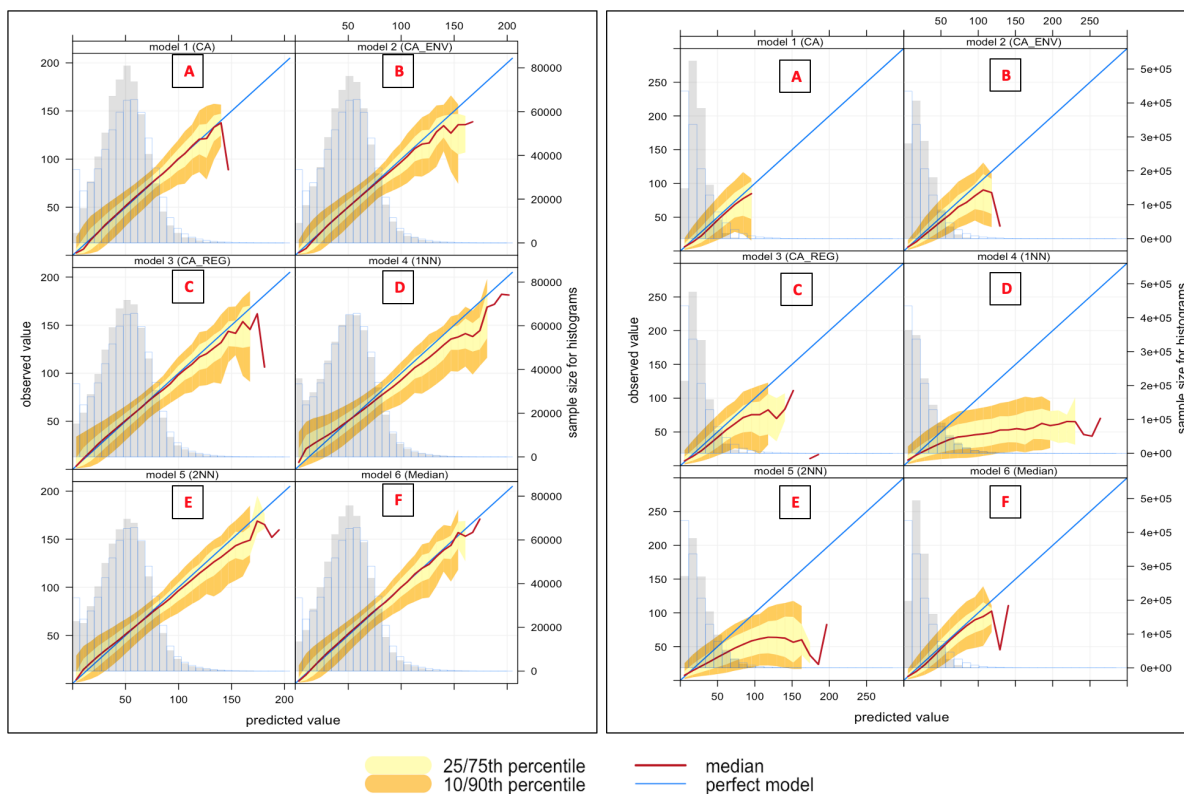
Imputation models	N	FAC2	MB	NMB	RMSE	R	IOA
<b>O<sub>3</sub></b>							
model 1 (CA)	71	0.867	-0.008	0	15.267	0.794	0.712
model 2 (CA+ENV)	71	0.877	1.113	0.022	14.627	0.815	0.729
model 3 (CA+REG)	71	0.872	-0.011	0	15.014	0.807	0.723
model 4 (INN)	71	0.831	-1.179	-0.024	17.494	0.757	0.681
model 5 (2NN)	71	0.871	-0.835	-0.017	15.159	0.808	0.721
model 6 (Median)	71	<b>0.888</b>	-0.373	-0.008	<b>13.776</b>	<b>0.837</b>	<b>0.745</b>
<b>NO<sub>2</sub></b>							
model 1 (CA)	157	0.628	0.009	0	18.33	0.514	0.599
model 2 (CA+ENV)	157	<b>0.708</b>	0.247	0.01	<b>15.989</b>	<b>0.665</b>	<b>0.661</b>
model 3 (CA+REG)	157	0.63	0.171	0.007	18.364	0.527	0.6
model 4 (INN)	157	0.605	2.277	0.095	22.591	0.464	0.533
model 5 (2NN)	157	0.618	2.774	0.116	20.46	0.494	0.558
model 6 (Median)	157	0.675	0.108	0.005	16.815	0.616	0.642
<b>PM<sub>2.5</sub></b>							
model 1 (CA)	77	0.835	-0.118	-0.012	5.265	0.787	0.713
model 2 (CA+ENV)	77	0.814	-0.064	-0.006	5.6	0.76	0.695
model 3 (CA+REG)	77	0.838	-0.064	-0.006	5.056	0.809	0.725
model 4 (INN)	77	0.791	0.058	0.006	5.536	0.79	0.7
model 5 (2NN)	77	0.823	0.02	0.002	4.952	0.823	0.726
model 6 (Median)	77	<b>0.854</b>	-0.144	-0.014	<b>4.745</b>	<b>0.831</b>	<b>0.743</b>
<b>PM<sub>10</sub></b>							
model 1 (CA)	75	0.86	-0.163	-0.01	8.747	0.668	0.667
model 2 (CA+ENV)	75	0.851	-0.148	-0.009	9.031	0.65	0.662
model 3 (CA+REG)	75	0.861	-0.043	-0.003	8.797	0.673	0.67
model 4 (INN)	75	0.816	0.113	0.007	10.363	0.608	0.627
model 5 (2NN)	75	0.858	0.106	0.006	9.23	0.661	0.668
model 6 (Median)	75	<b>0.882</b>	-0.216	-0.013	<b>8.224</b>	<b>0.715</b>	<b>0.697</b>



**Figure 1.** Taylor diagrams comparing modelled and observed concentrations for  $O_3$ ,  $NO_2$ ,  $PM_{2.5}$ , and  $PM_{10}$ .

These plots show how the modelled concentrations compare with the observed concentrations and how the models capture the variability in the concentrations. The spread of the modelled concentrations around the perfect model line (blue line) are shown by the shaded portions/quantile intervals. If narrow, it indicates high agreement/precision between the modelled and observed concentrations. The quantile intervals also represent the uncertainty bands. In some cases these intervals do not extend along with the median line due to insufficient concentrations to calculate them. The model with good performance is obtained when the median (red line) coincides with the perfect model (blue line) and when the spread in the percentile is as narrow as possible.

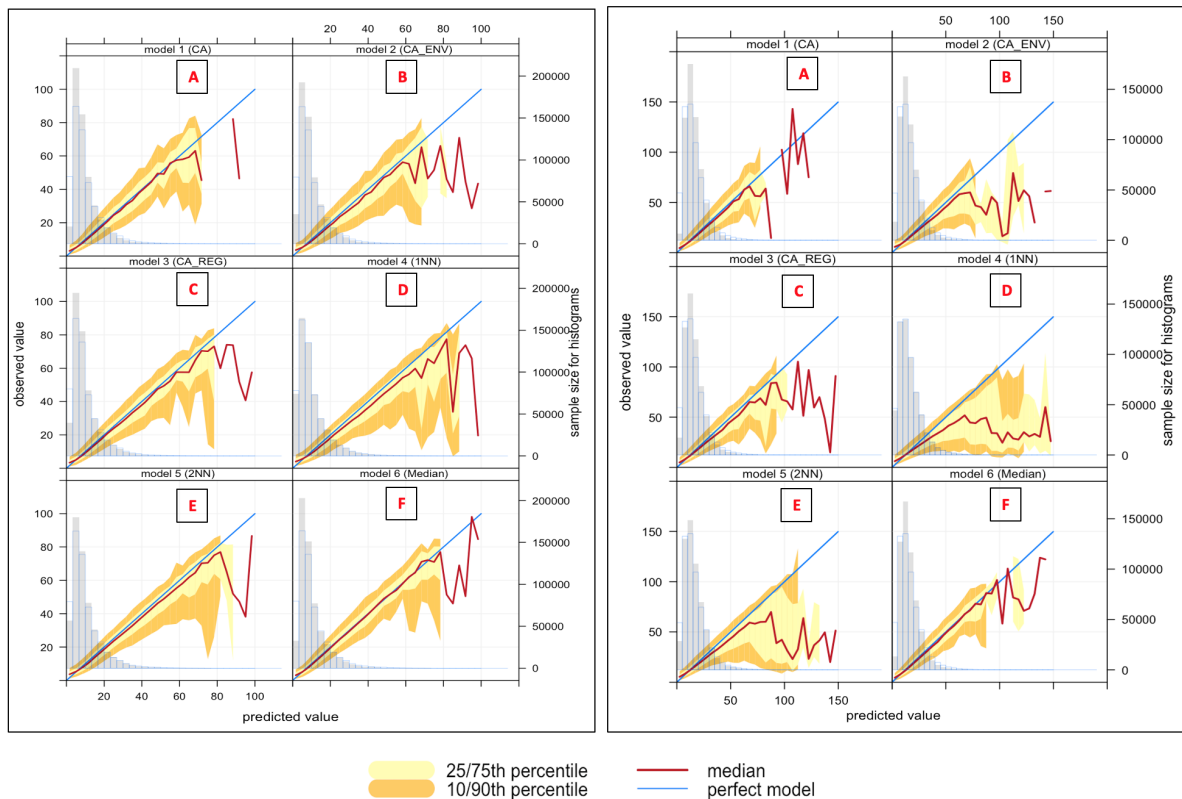
From these plots, in general, the histograms indicate that model 4 (1NN) (Panel D) has better estimation of the variability between the observed and modelled concentrations, as observed before, even though the median line does not match the perfect model. This model is positively biased at high concentrations, as shown by the departure of the median line below the blue line for all pollutants. This result supports our analysis from the Taylor diagram that model 4 (1NN) has the lowest variability between modelled and observed concentrations, but with lower correlation coefficient and the highest-centred root means squared for all pollutants.



**Figure 2.** Conditional quantile plot of modelled and observed pollutants concentrations of O<sub>3</sub> (left plot) and NO<sub>2</sub> (Right) for proposed imputation models; (A) model 1 (CA), (B) model 2 (CA+ENV), (C) model 3 (CA+REG), (D) model 4 (1NN), (E) model 5 (2NN), (F) model 6 (Median).

In Fig. 2, (left plot) the O<sub>3</sub> models show that most modelled concentrations match the observations well for a wide range of values. The histograms indicate underestimation in general at the extreme low and high concentrations. In general, the cluster and median imputation methods (i.e. that use averaging) will tend to struggle to reproduce the lowest and highest concentrations since they take an average approach. Moreover, the highest concentrations are typically limited to relative few data points. The cases of the high ozone concentrations typically occur during specific meteorological conditions and are episodic in nature, and there may be small differences in timings of the peak concentrations at different sites. The very low ozone concentrations are likely to occur at specific sites (near to roads where emissions of nitric oxide are large) and so may not be reproduced in the models which take a cluster average or where a nearest neighbour site is not a similar type of site.

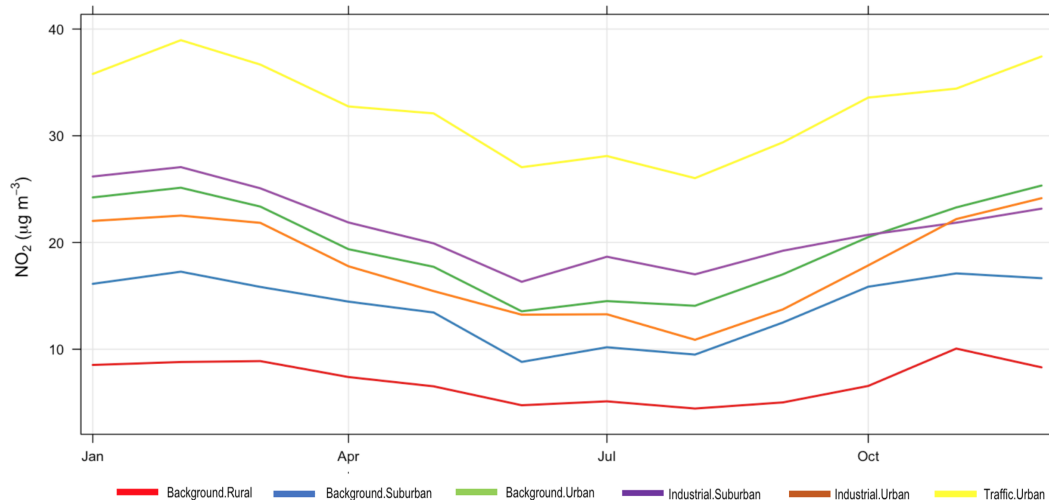
Model 6 (Median) (Panel F) has the best performance indicated by an overlapping median line with the blue line. This model has the lowest mean bias and the highest degree of agreement indicated by the narrow spread of the modelled concentration quantile intervals.



**Figure 3.** Conditional quantile plot of modelled and observed pollutants concentrations of PM<sub>2.5</sub> (left plot) and PM<sub>10</sub> (right plot) for proposed imputation models; (A) model 1 (CA), (B) model 2 (CA+ENV), (C) model 3 (CA+REG), (D) model 4 (1NN), (E) model 5 (2NN), (F) model 6 (Median).

In the same figure (right plot), NO<sub>2</sub> models show different behaviours from this analysis. Even though the statistical analysis shows that model 2 (CA+ENV)(Panel B) gives the best performance, it is clear that in this model, the modelled concentrations tend to be lower than observations for most concentration levels (the medians are under the blue line) and the width of the 10/75th and 10/90th percentiles is quite broad. The only advantage of using this model is its ability to capture a wide range of concentrations. Model 4 (1NN) (Panel D) compared to other models can reproduce the higher concentrations (higher than 125  $\mu\text{g m}^{-3}$ ) as it does not take an average approach. However, this model is positively biased (NMB = 0.095), which is shown by the departure of the median line from the blue one.

The variation between PM<sub>2.5</sub> models in Fig.3 (left plot) show similar performance for the different models. The quantile intervals are wider within the area of high concentrations  $\leq 60\mu\text{g m}^{-3}$ , and all models underestimate the high concentrations  $\leq 80\mu\text{g m}^{-3}$ , note that these concentrations are very low frequency events.



**Figure 4.** Monthly average concentrations of observed NO<sub>2</sub> for each environmental types for year 2018.

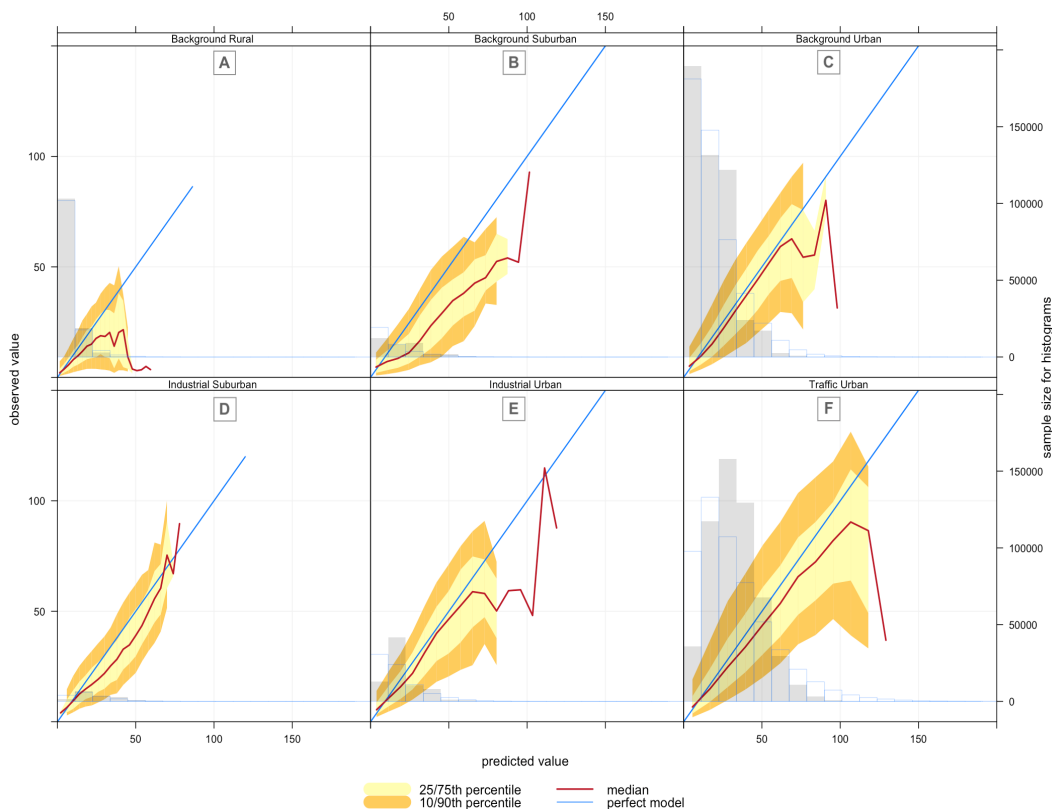
Model 6 (Median) (Panel F) gives better performance indicated by the narrow spread of the modelled concentration quantile intervals and minimal bias, indicated by the overlaps between the red and blue lines compared to other models. Models for  
 280 PM<sub>10</sub> (right plot) show similar performance to PM<sub>2.5</sub>.

#### 5.1.4 Model evaluation based on conditional quantile analysis and station environmental types

In this analysis, we focus on the performance of model 6 (Median) and model 2 (CA+ENV), as those performed best for the different pollutants in the previous section, but now we break down the analysis for the six environmental types (background rural, background urban, background suburban, and industrial urban, industrial suburban, and traffic urban) to which stations  
 285 belong. Notice that a pollutant may/may not be measured in all stations and the number of stations of each type is different as shown in Table 2. We also use conditional quantiles to analyse our model's performance within each environmental type.

First, we show the monthly average concentrations for each pollutant under each environment type in our test dataset (year 2018), to help understand the normal variation of the pollutant concentrations in different environment types. Figs. 5, 7, 9, and 11 show conditional quantile plots by the environmental types for the selected models. Table 2 shows the statistical measures  
 290 of performance also broken down by environment type.

The most common sources for NO<sub>2</sub> are roads, however NO<sub>2</sub> concentrations are influenced by traffic density, road locations, and meteorological conditions, which cause variation from one roadside location to another. Fig. 4 shows that high NO<sub>2</sub> concentrations are found at traffic urban followed by industrial suburban, then background urban sites, while the background rural sites have the lowest NO<sub>2</sub> concentrations.



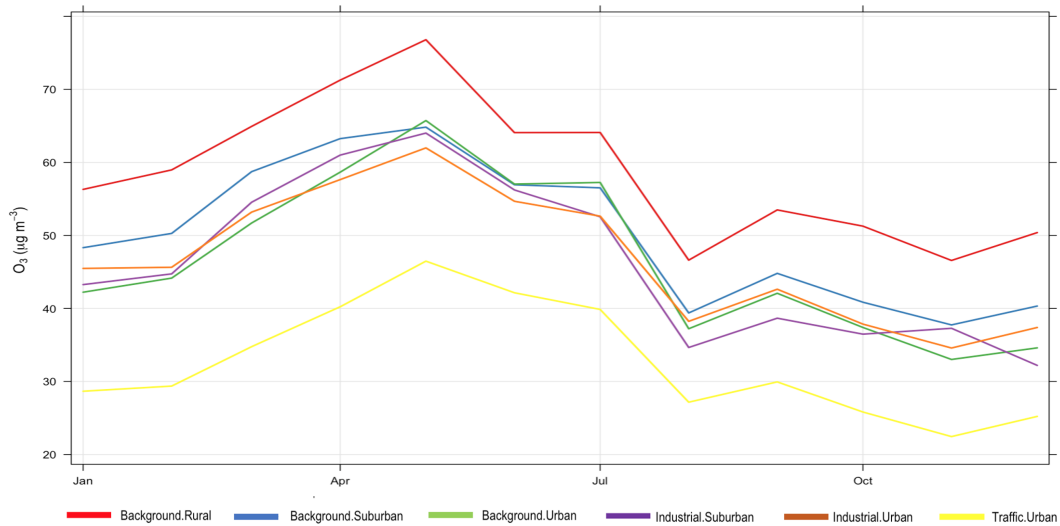
**Figure 5.** Conditional quantile plot of modelled and observed pollutants concentrations of NO<sub>2</sub> based on model 2 (CA+ENV) for all station environmental types, (A) background rural, (B) background suburban, (C) background urban, (D) industrial suburban, (E) industrial urban and (F) traffic urban.

295 Fig.5 shows the conditional quantile plots by station type for NO<sub>2</sub> imputation using model 2 (CA+ENV). Here, we see that modelled concentrations are higher than observed concentrations with all environmental types. This is confirmed by all the statistical model quality measures presented in Table 2 where we can observe positive mean bias for NO<sub>2</sub>.

As NO<sub>2</sub> distributions in general are skewed to the lower values, and our selected model model 2 (CA+ENV) is based on the average concentrations, the model performs better with lower concentrations.

300 From Table 2 based on model RMSE, the model's best performance is associated with background rural stations, while the worst performance is shown for traffic urban stations. Contrasting this with quantile plots, Fig.5 (Panel A) shows that for background rural stations the histogram and the median line show better performance with lower concentrations (less than 30  $\mu\text{g m}^{-3}$ ). On the other hand, for traffic urban stations (Panel F), the quantile intervals are wider within the area of high concentrations (higher than 25  $\mu\text{g m}^{-3}$ ), and the modelled concentrations tend to be lower than observed concentrations.

305 For O<sub>3</sub>, Fig. 6 shows the monthly average concentrations of observed O<sub>3</sub> concentrations at each environment type. From that we can see that ozone in all environment types follow a similar trend. However, background rural stations have the



**Figure 6.** Monthly average concentrations of observed O<sub>3</sub> for each environmental types for year 2018.

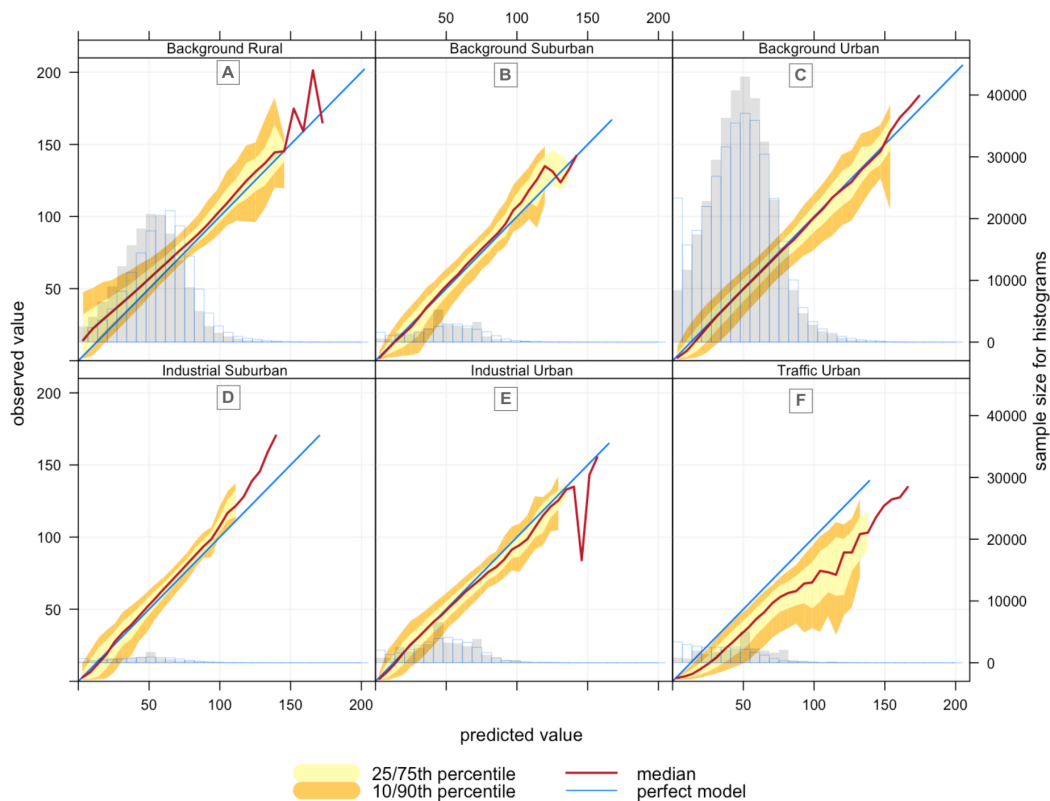
highest concentrations and traffic urban have the lowest, consistent with depletion of ozone due to rapid reaction with fresh emissions of nitric oxide from vehicles. Looking at model 6 (Median) performance in Table 2 based on the RMSE, the best model performance is associated with industrial urban stations and its average performance is associated with background rural stations (those with higher concentrations in Fig. 6), while its worst performance is associated with traffic urban stations (those with lower concentrations in Fig. 6).

Conditional quantile analysis in Fig. 7, shows the performance of model 6 (Median) for imputing O<sub>3</sub> for the six environmental types (Panels A to F). The model shows similar performance for industrial suburban (Panel D) and background rural stations (Panel A). For both types, the model is negatively biased (see also Table 2), meaning that the modelled concentrations tend to be lower than observed concentrations (the median lines are above the blue lines).

The worst performance based on the RMSE is associated with traffic urban stations (Panel F), which are the stations located at the roadsides. With those stations, the modelled concentrations are higher than observed concentrations, i.e. the modelled histogram is shifted to the right. This is indicated by the model positive bias (0.503). The median line also extends beyond the blue line, which means that some modelled concentrations are much higher than observed measurements.

The best model performance is associated with industrial urban stations (Panel E) according to the RSME, even though background urban stations (Panel C) appear to have the best performance by looking at the conditional quantile plots. The histogram of Panel C indicates that the distribution of the observed and modelled concentrations tend to be closer to each other for higher concentrations. However, the model overestimates the average concentrations at these stations (between 25 to 70 µg m<sup>-3</sup>) and underestimates the very low concentrations.

Fig. 8, shows PM<sub>2.5</sub> concentrations at rural areas are lower than those at suburban, urban background and traffic urban areas. That is consistent with the model performance at these sites. Fig. 9 shows corresponding conditional quantile plots by station

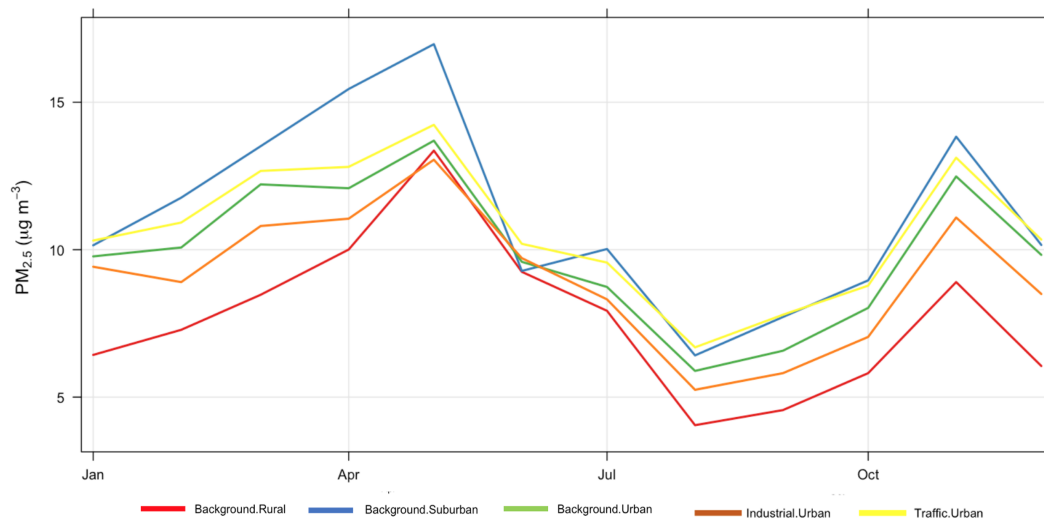


**Figure 7.** Conditional quantile plot of modelled and observed pollutants concentrations of O<sub>3</sub> based on model 6 (Median) for station environmental types, (A) background rural, (B) background suburban, (C) background urban, (D) industrial suburban, (E) industrial urban and (F) traffic urban.

types. Imputing PM<sub>2.5</sub> concentrations using model 6 (Median) gives similar performance for the different station types. In general, the model underestimates the concentrations of PM<sub>2.5</sub> especially for high concentration levels. Table 2, shows that the model underestimates high concentrations at suburban, urban background and traffic urban areas, indicated by the model  
 330 negative biases, while it overestimates the concentrations at industrial urban and background rural sites. The model shows worst performance for traffic urban (panel E), and this is also indicated by the highest RMSE (5.098) shown in Table 2. The model underestimates the concentrations at these stations, which is confirmed by the model bias (-0.073) in Table 2. On the other hand, the model's best performance is associated with background suburban sites (Fig. 9 (panel B)), even though it underestimates PM<sub>2.5</sub> concentrations with a mean bias of (-0.013).

335 Finally, PM<sub>10</sub> levels at background rural and urban areas are lower than those at industrial and traffic urban areas as shown in Fig. 10. For PM<sub>10</sub>, imputation performance shown in Fig.11 is similar for background urban and background rural sites (panels A and B). The model overestimates the concentrations of PM<sub>10</sub> that are  $\leq 10\mu g m^{-3}$ , while it underestimates the high





**Figure 8.** Monthly average concentrations of observed PM<sub>2.5</sub> for each environmental types for year 2018.

concentrations of PM<sub>10</sub> at industrial urban (slightly) and traffic urban sites (panels C and D). That is confirmed by the model mean bias at these sites (-0.002, -0.106 ) as shown on Table 2.

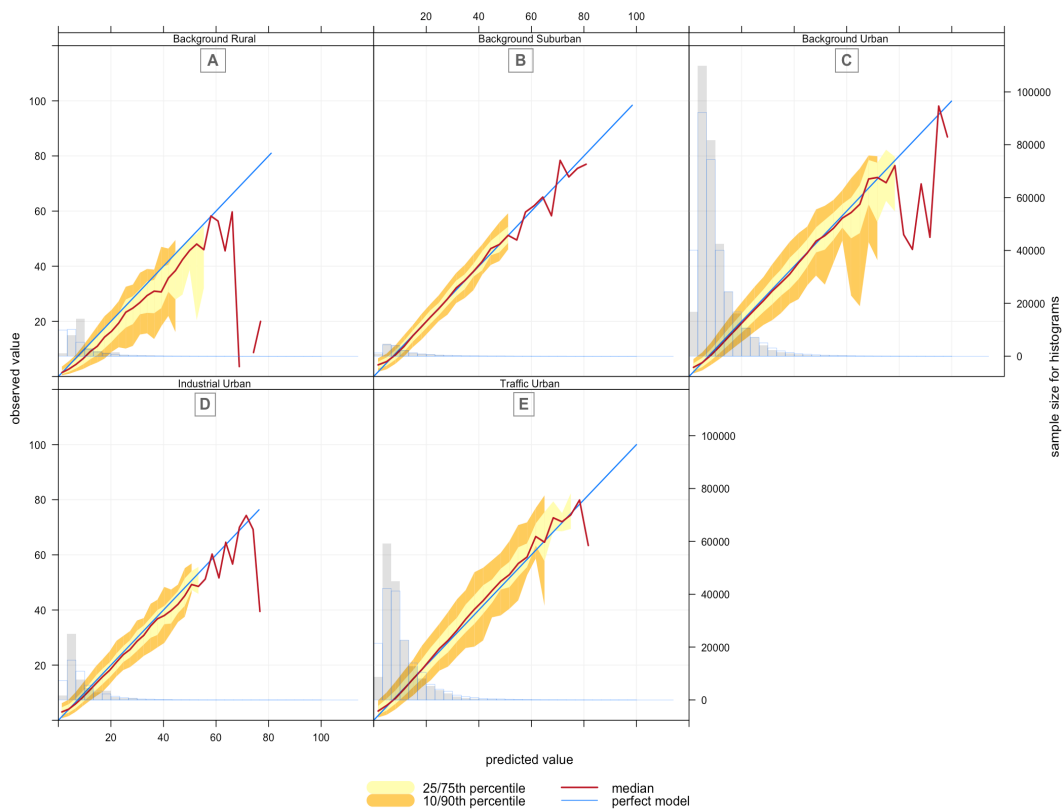
340 Next, we show some examples of our imputed TS compared to the real TS for each pollutant using the selected imputation models in some stations. The following examples in Figs. 12, 13, 14, and 15 show the observed and imputed hourly pollutants concentrations for the four pollutants using the selected imputation model that gave better imputation. We also apply the models where there is a period of missing values in the observed concentrations to give some idea of how the models work for the whole TS including when real values do not exist.

345 Fig. 12, compares the observed hourly concentrations of PM<sub>2.5</sub> (red) at 'London Eltham' station for 1/1/2018 to 15/1/2018 with imputed concentrations (black) using model 6 (median). As we can see, the variation between the imputed and the real TS is very small and the imputed TS reproduces the trend very well, even though, there is a period of missing concentrations within the observed TS (red). Similarly Fig. 13 represents the observed hourly concentrations of PM<sub>10</sub> (red) at 'Oxford St Ebbs' station for the same period of time with imputed concentrations using model 6. We can see that model 6 underestimates the high concentrations and overestimates the very low concentrations of PM<sub>10</sub> and PM<sub>2.5</sub>, as mentioned previously in the analysis in Sec. 5.1.3. However, there is still a good match of the trend.

350

Fig. 14 shows a comparison of imputed (black) and observed (red) TS for NO<sub>2</sub> concentrations at 'Birmingham Acocks Green' station for the same period of time (1/1/2018 to 15/1/2018), but produced by a different imputation model (model 3 (CA+ENV)) that gives better imputation than others for NO<sub>2</sub>. It is known that NO<sub>2</sub> has greater spatial variability than other pollutants and it is very complex to impute, the variation between the imputed and the real TS is slightly higher when compared to the previous examples.

355



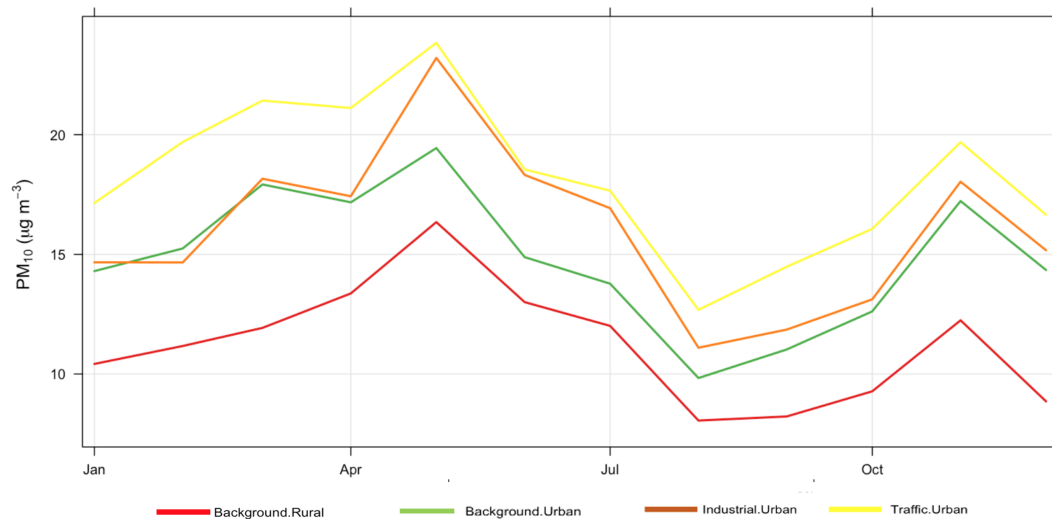
**Figure 9.** Conditional quantile plot of modelled and observed pollutants concentrations of  $PM_{2.5}$  based on model 6 (Median) for station environmental types, (A) background rural, (B) background suburban, (C) background urban, (D) industrial suburban, (E) traffic urban.

Fig.15 shows a comparison of the imputed (black) and observed (red) TS for  $O_3$  concentrations at 'Birmingham Acocks Green' station for the period (16/1/2018 to 23/1/2018), produced by model 6 (median). The imputation underestimate the concentrations but represents the trends of high and low values.

## 360 5.2 Evaluating the imputed concentrations based on the Daily Air Quality Index (DAQI)

After imputing the measured pollutants in all the stations, we calculate DAQI from the imputed data, as explained in Section. 4.3. Then we compare it with the DAQI from the observed data to see our selected models' performances. The selected models are model 6 (Median) for  $O_3$ ,  $PM_{2.5}$ , and  $PM_{10}$  and Model 2 (CA+ENV) for  $NO_2$ .

We compare the imputed DAQI with the observed DAQI based on RMSE, and the number of days where there are agreements and disagreements. The total number of days in our data set is 60,955 days (167 stations \* 365 days), there are 2,212 days with missing observed DAQI (DAQI = 0) that have resulted from missing observations on those days. The total number of days to compare is 58,743 days.



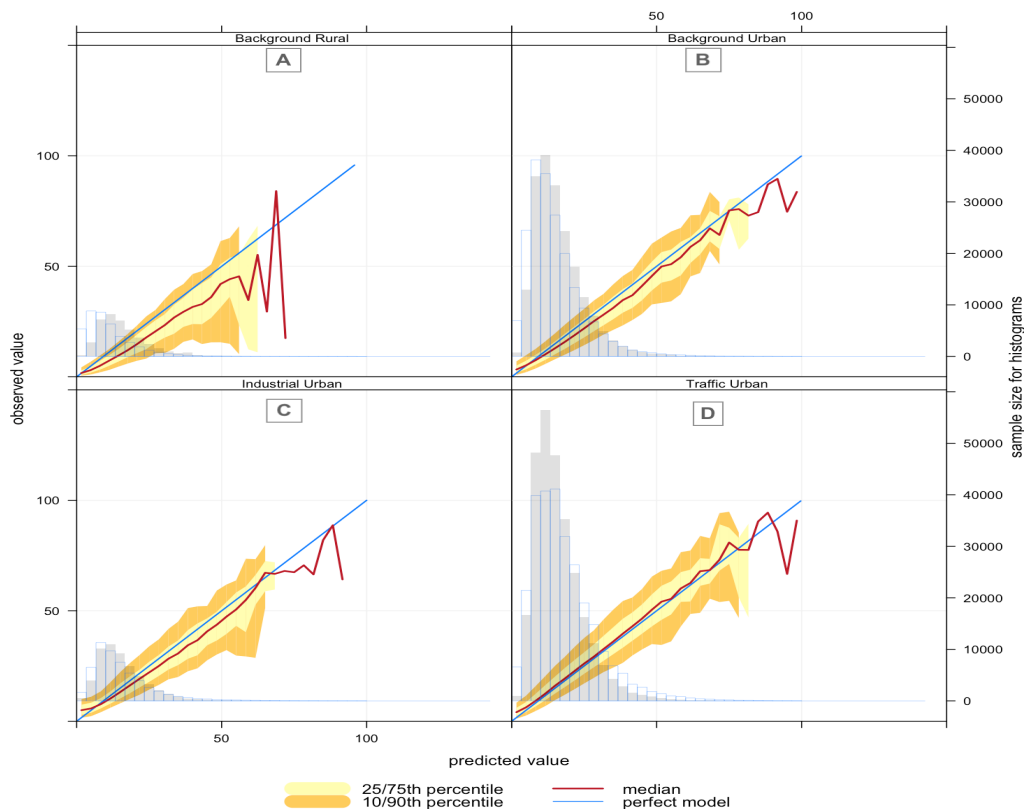
**Figure 10.** Monthly average concentrations of observed PM<sub>10</sub> for each environmental types for year 2018.

In general, the total average of RMSE from all days in all stations is (0.55). As the station type and the region may affect our imputation, Fig. 16 shows the average RMSE based on air quality regions in the UK (Panel A), and station environmental types (Panel B); the size of the circles are the number of stations at each type. Panel A shows that stations classed as Traffic Urban are associated with the highest RMSE (0.62), while Industrial Suburban stations have the lowest RMSE (0.36). Panel B shows that the North East region is associated with the lowest RMSE (0.44), while South Wales has the highest RMSE (0.74) between imputed and observed DAQI.

We also study the correlation between number of measured pollutants in a station and the agreement between modelled and observed DAQI to see if number of measured pollutants impacts our model's performance.

First, we classify stations based on number of measured pollutants to: stations that measured one, two, three and all four pollutants, as shown in Table 3. Each row in this table represents one group. The second column is the total number of days with associated DAQI from all stations in each group. The RMSE and index of agreement (IOA) are the average of errors and the degree of agreement between observed and modelled DAQI from all stations in each group, then the percentage of each pollutant in each group. Based on this table, we find that stations that measure four pollutants have the lowest RMSE (0.506) and the highest (IOA) (0.806), while stations that measured one pollutant have the worst performance. The majority of stations with one pollutant are stations that measure NO<sub>2</sub> with (87%) of total number of station in this group (50 stations).

We also compare the imputed and the observed DAQI based on the number of days where the imputed DAQI agrees and disagrees with the observed DAQI. Table. 4, shows those results and the percentage of time that these situations occurred, means when agreement/disagreement is found for each DAQI. The total number of days where the imputed DAQI agrees with the observed DAQI is 43,906 day (75%), while there is 14,837 (25%) days of disagreement. We classify the disagreement into two types: the imputed DAQI is higher or lower than the observed DAQI. We find that there are 10,916 days, where



**Figure 11.** Conditional quantile plot of modelled and observed pollutants concentrations of  $PM_{10}$  based on model 6 (Median) for station environmental types, (A) background rural, (B) background urban, (C) industrial urban, (D) traffic urban.

the imputed DAQI is lower than the observed DAQI, and 3,921 days, where the imputed DAQI is higher than the observed DAQI. In most cases, the imputed DAQI is lower than the observed DAQI, in accordance with our analysis of the imputation  
 390 models that showed underestimation of the pollutant concentrations. From this table, we can see that the highest percentage of disagreement is 42.96% of total number of disagreement (14837) when observed DAQI is 2 and imputed DAQI is 1, followed by 21.35% of disagreement when observed DAQI is 3 and imputed DAQI is 2.

## 6 Discussion and Conclusions

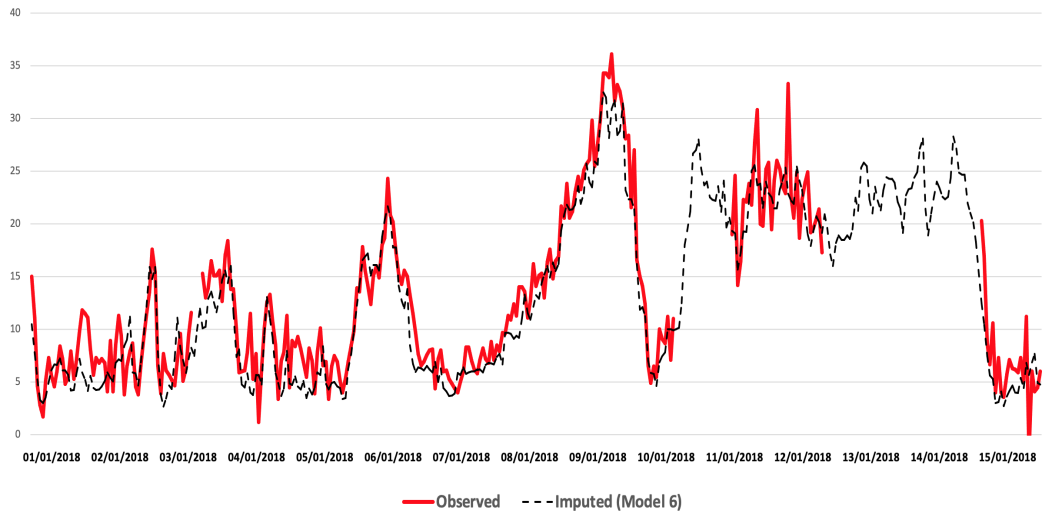
In this work, we evaluated our proposed models to impute missing pollutants in a station based on statistical and graphical  
 395 model evaluation functions (Taylor's diagrams and Conditional quantile plots), that are designed to evaluate air pollution modelling. We found that the best imputation model based on statistical analysis is model 6 (Median) for  $O_3$ ,  $PM_{10}$ , and  $PM_{2.5}$  and Model 2 (CA+ENV) for  $NO_2$  imputation. The station environmental type plays an essential role with  $NO_2$  imputation, because  $NO_2$  shows local patterns, as it is concentrated where it is emitted in urban areas and near to the roadside. Adding

**Table 2.** Performance of the hourly pollutant concentrations imputation models using model 6 (Median) for O<sub>3</sub>, PM<sub>2.5</sub>, and PM<sub>10</sub> and Model 2 (CA+ENV) for NO<sub>2</sub> based on statistical measures for all station environment types for all pollutants.

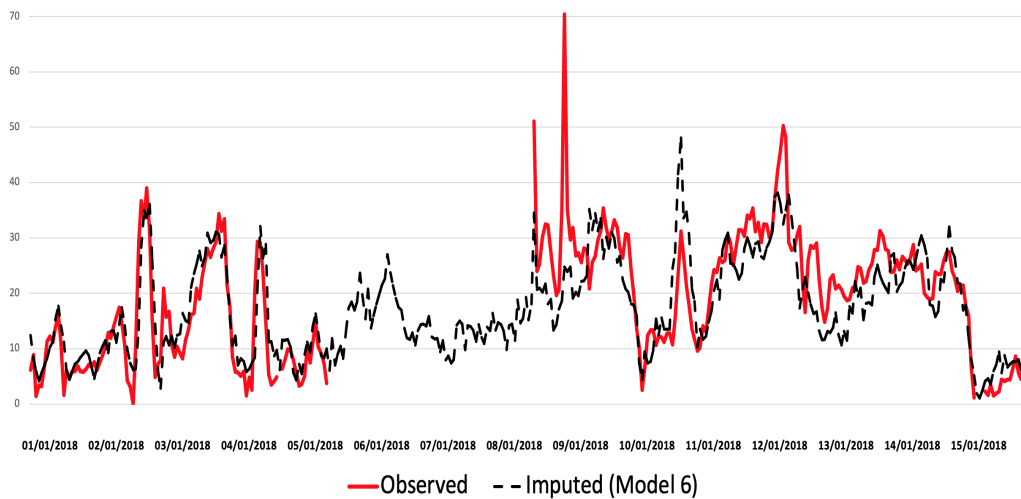
Imputation models	Environment Type	N	MB	NMB	RMSE
O <sub>3</sub>					
model 6 (Median)	Background Rural	19	-7.648	-0.130	14.967
model 6 (Median)	Background Suburban	3	0.133	0.003	12.530
model 6 (Median)	Background Urban	39	1.578	0.034	12.780
model 6 (Median)	Industrial Suburban	2	-1.392	-0.030	11.311
model 6 (Median)	Industrial Urban	4	1.561	0.033	11.273
model 6 (Median)	Traffic Urban	3	16.456	0.503	21.580
NO <sub>2</sub>					
model 2 (CA+ENV)	Background Rural	15	0.060	0.008	6.699
model 2 (CA+ENV)	Background Suburban	5	6.590	0.470	13.576
model 2 (CA+ENV)	Background Urban	58	0.025	0.001	12.442
model 2 (CA+ENV)	Industrial Suburban	4	3.929	0.181	11.939
model 2 (CA+ENV)	Industrial Urban	11	0.235	0.013	10.481
model 2 (CA+ENV)	Traffic Urban	65	-0.014	0.000	20.500
PM <sub>2.5</sub>					
model 6 (Median)	Background Rural	5	2.167	0.292	5.004
model 6 (Median)	Background Suburban	2	-0.143	-0.013	3.434
model 6 (Median)	Background Urban	41	-0.072	-0.007	4.685
model 6 (Median)	Industrial Urban	6	0.080	0.009	3.982
model 6 (Median)	Traffic Urban	23	-0.781	-0.073	5.098
PM <sub>10</sub>					
model 6 (Median)	Background Rural	5	4.205	0.369	8.036
model 6 (Median)	Background Urban	26	1.236	0.082	7.097
model 6 (Median)	Industrial Urban	7	-0.037	-0.002	10.027
model 6 (Median)	Traffic Urban	37	-1.939	-0.106	8.586

to that NO<sub>2</sub> is shorter lived than other pollutants and shows greater spatial variability, with concentrations being strongly influenced by the environment type (e.g. roadside, urban background, rural). This changes the NO<sub>2</sub> concentrations from one location to another based on the environmental type (CenterForCities, 2020).

On the other hand, the graphical model evaluation functions showed these models' performance based on the distribution of the concentrations and the degree of agreement between imputed/modelled and observed concentrations. These functions

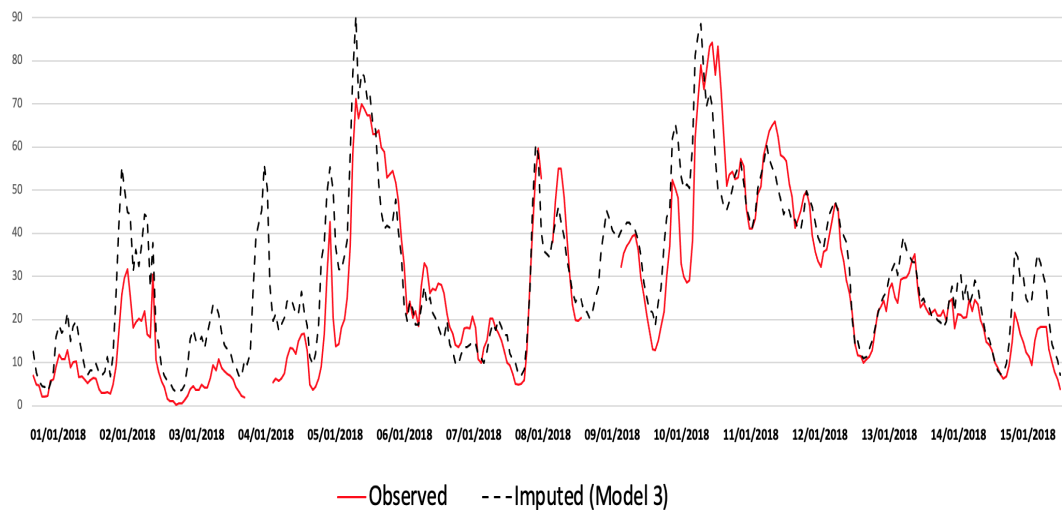


**Figure 12.** Imputed (black) and real (red) TS comparison for PM<sub>2.5</sub> at 'London Eltham' station from 1/1/2018 to 15/1/2018.

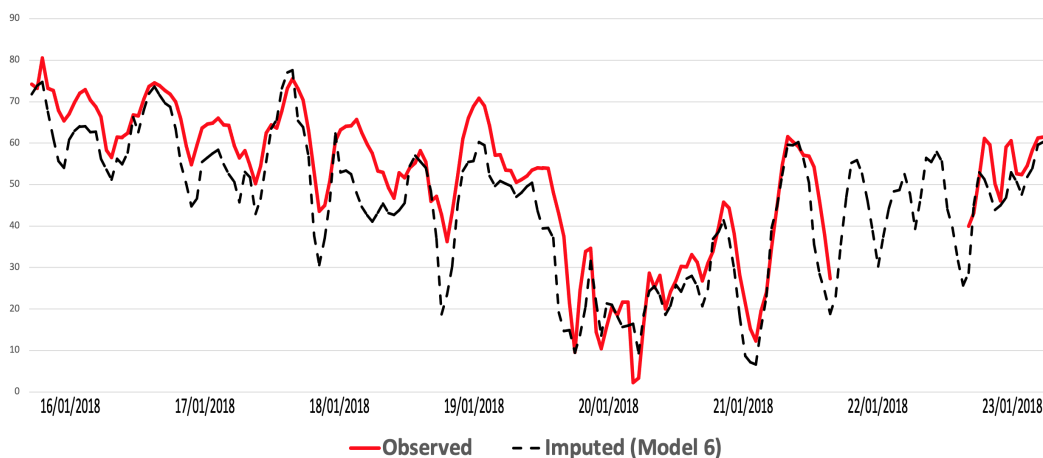


**Figure 13.** Imputed (black) and real (red) TS comparison for PM<sub>10</sub> at 'Oxford St Ebbes' station from 1/1/2018 to 15/1/2018.

405 help us to understand the relationship between the distributions of the observations and the model's performance. From the histograms in Figs. 2 and 3 we noted that the overall distribution of the observed concentrations of NO<sub>2</sub>, PM<sub>2.5</sub>, and PM<sub>10</sub> are skewed to the lower values, while O<sub>3</sub> has a more normal distribution. From these histograms, we also noticed that the distributions of the modelled O<sub>3</sub> concentrations are shifted to lower values, while other pollutants modelled concentrations are shifted to higher values. Hence the model is not always able to reproduce the edges of the distribution correctly. The skewness



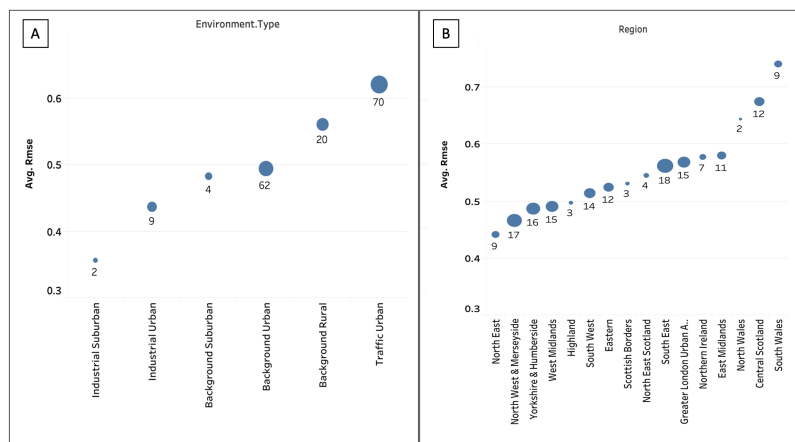
**Figure 14.** Imputed (black) and real (red) TS comparison for NO<sub>2</sub> at 'Birmingham Acocks Green' station from 1/1/2018 to 15/1/2018.



**Figure 15.** Imputed (black) and real (red) TS comparison for O<sub>3</sub> at 'Birmingham Acocks Green' station from 16/1/2018 to 23/1/2018.

in modelled values is mostly associated with model 1 (CA). As a consequence, this shows the greatest difference in skewness  
 410 between the distributions of the observed and modelled values. However, model 1 (CA) in combination with others, as part of  
 model 6 (Median), reduces the skewness in modelled values and generates better imputation, resulting in the lowest RMSE.

Model 6 (Median) is based on the median concentrations from stations with temporal and spatial similarity, so this model's  
 expected performance is to underestimate the highest values and overestimate the lowest values with a normal distributed



**Figure 16.** The model performance based on DAQI RMSE: (A) the average of the RMSE based on station environmental types, (B) the average of the RMSE based on air quality regions.

**Table 3.** Comparing observed and modelled DAQI based on number of measured pollutants in stations.

Number of measured pollutants	Number of days in all stations	Number of Stations	RMSE	IOA	Percentage (O <sub>3</sub> )	Percentage (NO <sub>2</sub> )	Percentage (PM <sub>2.5</sub> )	Percentage (PM <sub>10</sub> )
1	15684	50	0.542	0.769	6.1	87.8	2.0	4.1
2	17581	48	0.583	0.756	24.5	48.0	8.2	19.4
3	15443	43	0.516	0.814	14.0	31.8	32.6	21.7
4	9398	26	0.496	0.814	25.0	25.0	25.0	25.0

dataset. We found that the models performance can vary based on the environmental type and the nature of the pollutant, as shown in our analysis of model performance and DAQI RMSE.

Through our analysis, we also found that the variation of the model’s performance with different environmental types is due to the pollutant behaviour and its emitted sources.

Model 6 (Median) performance with O<sub>3</sub> imputation changes from one environmental type to another due to the ozone’s behaviour at these locations. As we know, ozone is not directly emitted into the air, but it is formed as a secondary pollutant by chemistry involving nitrogen oxides (NO<sub>x</sub>), the sum of NO<sub>2</sub>, nitric oxide (NO) and volatile organic compounds (VOCs) in the presence of sunlight (Diaz et al., 2020). This chemistry is non-linear and newly emitted NO can react with O<sub>3</sub> leading to reductions in O<sub>3</sub> concentrations close to sources of NO (e.g. in urban areas and in particular, close to roads). Consequently, ozone concentrations in urban areas are often lower than those at rural areas (H. Khan et al., 2017), as shown in Fig. 6.

Fig. 7 (Panel A), shows see that the model produces a distribution shifted to the left toward lower values, not capturing the ozone for rural areas that are associated with higher concentrations of O<sub>3</sub>. Similarly, industrial suburban stations (Panel D) have a higher frequency of high concentrations (higher than 25 μg m<sup>-3</sup>), as shown in the histogram (Panel D). Note that



**Table 4.** Number of days where imputed DAQI agrees/disagrees with observed DAQI.

Index Agreement							
Observed DAQI	Imputed DAQI	Number of days	Percentage of Days	Observed DAQI	Imputed DAQI	Number of days	Percentage of Days
1	1	18920	43.09	5	5	110	0.25
2	2	16351	37.24	6	6	19	0.04
3	3	7969	18.15	7	7	5	0.01
4	4	525	1.20	8	8	7	0.02
<b>Total agreement</b>							43906
<b>Total Percentage</b>							74.743
Index Disagreement							
Observed DAQI	Imputed DAQI	Number of days	Percentage of Days	Observed DAQI	Imputed DAQI	Number of days	Percentage of Days
1	2	1818	12.25	5	2	6	0.04
1	3	255	1.72	5	3	54	0.36
1	4	11	0.07	5	4	203	1.37
1	5	2	0.01	5	6	12	0.08
1	8	1	0.01	6	7	4	0.03
2	8	1	0.01	6	2	2	0.01
2	1	6374	42.96	6	3	5	0.03
2	3	1479	9.97	6	4	31	0.21
2	4	10	0.07	6	5	45	0.30
2	5	4	0.03	7	8	1	0.01
3	1	337	2.27	7	2	1	0.01
3	2	3168	21.35	7	3	2	0.01
3	4	241	1.62	7	4	1	0.01
3	5	18	0.12	7	5	10	0.07
3	6	2	0.01	7	6	11	0.07
4	1	17	0.11	8	7	3	0.02
4	2	38	0.26	8	9	1	0.01
4	3	598	4.03	8	3	2	0.01
4	5	58	0.39	8	6	1	0.01
4	6	2	0.01	9	8	2	0.01
5	7	1	0.01	10	7	2	0.01
5	1	2	0.01	10	8	1	0.01
<b>Total disagreement</b>							14837
<b>Total Percentage</b>							25.257

majority of stations measuring O<sub>3</sub> are background rural or background urban, with few stations in other categories. While with traffic urban (Panel F), where the model performs the worst, some modelled concentrations are much higher than observed measurements. This lack of fit may be explained because ozone is suppressed by new emissions of NO close to sources (traffic) which reduce the amount of O<sub>3</sub> in those station types.

From the same figure (Panel C), as shown in the histogram for background urban stations, there is a high frequency of low concentrations (less than 10 µg m<sup>-3</sup>) at these stations that the model does not capture. This is consistent with the reaction of newly emitted NO from urban roadside that reduces the concentrations of ozone at urban areas. Based on the RMSE and NMB, the model is a middle performing model. As shown in Table 2, the majority of stations measuring O<sub>3</sub> belong to this type.

As, we mentioned earlier that NO<sub>2</sub> is short lived so it has large differences between sites near sources (roadside) and those further away. Based on the RMSE Model 2 (CA+ENV) performs better with lower NO<sub>2</sub> concentrations than high values, and since these high NO<sub>2</sub> values exist near to traffic, the model performs the worst with traffic urban stations as shown in Fig.

5 (Panel F). In contrast, the model best performance is associated with background rural stations that have the lowest NO<sub>2</sub> concentrations.

440 PM<sub>2.5</sub> and PM<sub>10</sub> have many varied sources so in roads and industrial sites it can be associated with local sources, for example widespread primary sources (direct emissions) and diffused secondary sources (i.e. produced in the atmosphere following emissions of precursor gases). Whilst PM concentrations are often greater at roadside (DEFRA LAQM , 2021), the particles can have lifetimes of several days in the atmosphere, meaning that they can be distributed widely. The larger particles are subject to greater loss via sedimentation, so PM<sub>2.5</sub> is more evenly distributed than PM<sub>10</sub> (National Statistics, 2020). This behaviour  
445 can also be observed with Model 6 (Median) performance, where there are less variation with the model performance under different environment types compared to the variation of NO<sub>2</sub> and O<sub>3</sub>, as shown in Table 2.

We also observed that the distributions of NO<sub>2</sub>, PM<sub>2.5</sub> and PM<sub>10</sub> are skewed to lower concentrations which impact model performance at higher concentrations. All models perform worse for high concentrations with NO<sub>2</sub>, PM<sub>2.5</sub> and PM<sub>10</sub> than O<sub>3</sub>, indicated by the width of the quantiles at high values as shown in Figs. 2 and 3. Similarly, for lower concentrations, these  
450 models tend to perform better for NO<sub>2</sub>, PM<sub>2.5</sub>, and PM<sub>10</sub> than for O<sub>3</sub>. However, our selected models (model 6 (Median) and model 2 (CA+ENV)) are able to overcome this impact slightly.

Our approach enables us to impute/estimate plausible concentrations of multiple pollutants at stations across the UK, and the modelled concentrations from the selected models correlated well with the observed concentrations. The performance of these models is very good with a slight underestimation in model 6 (Median), especially with high concentrations. At the opposite  
455 end, Model 2 (CA+ENV) slightly overestimates the NO<sub>2</sub> concentrations, due to the regional behaviour of this pollutant.

We also analysed the performance of these models based on the daily modelled concentrations under different weather types using Lamb Weather Types (LWTs), which are a synoptic classification of daily weather patterns across the UK Lamb (1972). We found that these models work equally well for all LWTs, so we did not include this analysis in this work.

In conclusion, MVTS clustering enables imputation even when no measurement is available for a given pollutant since the  
460 station can be allocated to a cluster based on the value of the other pollutants measured. Our proposed imputation models, model 6 (Median) for O<sub>3</sub>, PM<sub>10</sub>, and PM<sub>2.5</sub> and Model 2 (CA+ENV), give the best performance for imputing these pollutants. The advantage of these model is that they aggregate the spatial and temporal imputation. The spatial imputation is obtained from the nearest stations and the temporal imputation is obtained by MVTS clustering that clusters the stations based on similarity in time.

465 In our future work, we aim to improve our imputation by considering more information about the stations, such as station altitude and location in relation to the weather effects. We may also consider the correlation between pollutants in our imputation, and include further analysis for the daily air quality index (DAQI), especially for those days when there is a variation between imputed and observed DAQI. Finally, we need to study all possible uncertainty associated with this type of application, since the pollution level may change from year to year due to some pollution episodes caused by high temperature, wind, wildfire or  
470 other reasons.

*Author contributions.* The experimentation and initial draft were produced by Wedad Alahamade as part of her PhD. Prof. Iain Lake and Prof. Claire Reeves contributed ideas, co-supervised the PhD and revised the draft manuscripts. Dr. Beatriz de la Iglesia was the main supervisor for the work and contributed to the draft revisions.

*Competing interests.* The authors declare that they have no conflict of interest.

## 475 **References**

- Alahamade, W., Lake, I., Reeves, C. E., and De La Iglesia, B.: Clustering Imputation for Air Pollution Data, in: International Conference on Hybrid Artificial Intelligence Systems, pp. 585–597, Springer, 2020.
- Alahamade, W., Lake, I., Reeves, C. E., and De La Iglesia, B.: A Multi-variate Time Series clustering approach based on Intermediate Fusion A case study in air pollution data imputation, *Neurocomputing*, accepted, 2021.
- 480 Austin, E., Coull, B. A., Zanobetti, A., and Koutrakis, P.: A framework to spatially cluster air pollution monitoring sites in US based on the PM<sub>2.5</sub> composition, *Environment international*, 59, 244–254, 2013.
- Carbajal-Hernández, J. J., Sánchez-Fernández, L. P., Carrasco-Ochoa, J. A., and Martínez-Trinidad, J. F.: Assessment and prediction of air quality using fuzzy logic and autoregressive models, *Atmospheric Environment*, 60, 37–50, 2012.
- Carlaw, D. C. and Ropkins, K.: Openair—an R package for air quality data analysis, *Environmental Modelling & Software*, 27, 52–61, 485 2012.
- CenterForCities(2020): Cities Outlook 2020 - Air quality in UK cities, <https://www.centreforcities.org/publication/cities-outlook-2020/>, 2020.
- DEFRA (2013): Daily Air Quality Index implementation Report, [https://uk-air.defra.gov.uk/library/reports?report\\_id=750/](https://uk-air.defra.gov.uk/library/reports?report_id=750/), 2013.
- DEFRA (2021): About Air Pollution, <https://uk-air.defra.gov.uk/air-pollution>, 2021.
- 490 DEFRA LAQM (2021): Public Health Sources and Effects of PM<sub>2.5</sub>, <https://laqm.defra.gov.uk/public-health/pm25.html>, 2016.
- Di Bello, G., Lapenna, V., Macchiato, M., Satriano, C., Serio, C., Tramutoli, V., et al.: Parametric time series analysis of geoelectrical signals: an application to earthquake forecasting in Southern Italy, 1996.
- Diaz, F. M., Khan, M. A. H., Shallcross, B., Shallcross, E. D., Vogt, U., and Shallcross, D. E.: Ozone Trends in the United Kingdom over the Last 30 Years, *Atmosphere*, 11, 534, 2020.
- 495 Dick Derwent, Andrea Fraser, J. A. M. J. P. W. and Murrells, T.: Report: Evaluating the Performance of Air Quality Models, Department for Environment, Food and Rural Affairs, London, 2010.
- Du, S., Li, T., Yang, Y., and Horng, S.-J.: Multivariate time series forecasting via attention-based encoder–decoder framework, *Neurocomputing*, 388, 269–279, 2020.
- D’Urso, P., De Giovanni, L., and Massari, R.: Robust fuzzy clustering of multivariate time trajectories, *International Journal of Approximate* 500 *Reasoning*, 99, 12–38, 2018.
- Fontes, C. H. and Budman, H.: A hybrid clustering approach for multivariate time series—a case study applied to failure analysis in a gas turbine, *ISA transactions*, 71, 513–529, 2017.
- H. Khan, M. A., Morris, W. C., Galloway, M., A. Shallcross, B. M., Percival, C. J., and Shallcross, D. E.: An Estimation of the Levels of Stabilized Criegee Intermediates in the UK Urban and Rural Atmosphere Using the Steady-State Approximation and the Potential Effects 505 of These Intermediates on Tropospheric Oxidation Cycles, *International journal of chemical kinetics*, 49, 611–621, 2017.
- Ignaccolo, R., Ghigo, S., and Giovenali, E.: Analysis of air quality monitoring networks by functional clustering, *Environmetrics*, 19, 672–686, 2008.
- Lam, N. S.-N.: Spatial interpolation methods: a review, *The American Cartographer*, 10, 129–150, 1983.
- Lamb, H. H.: British Isles weather types and a register of the daily sequence of circulation patterns 1861–1971, 1972.
- 510 Liao, T. W.: Clustering of time series data—a survey, *Pattern recognition*, 38, 1857–1874, 2005.

- National Statistics(2020): National Statistics Concentrations of Particulate Matter PM<sub>10</sub> and PM<sub>25</sub>, <https://www.gov.uk/government/publications/air-quality-statistics/concentrations-of-particulate-matter-pm10-and-pm25>, 2020.
- Sarda-Espinosa, A., Sarda, M. A., LazyData, T., Collate'DBA, R., and R'SBD, R.: Package 'dtwclust', 2017.
- 515 Seto, S., Zhang, W., and Zhou, Y.: Multivariate time series classification using dynamic time warping template selection for human activity recognition, in: 2015 IEEE Symposium Series on Computational Intelligence, pp. 1399–1406, IEEE, 2015.
- Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *Journal of Geophysical Research: Atmospheres*, 106, 7183–7192, 2001.
- Tuysuzoglu, G., Birant, D., and Pala, A.: Majority Voting Based Multi-Task Clustering of Air Quality Monitoring Network in Turkey, *Applied Sciences*, 9, 1610, 2019.
- 520 Wickham, H. and Wickham, M. H.: Package tidyverse, Easily Install and Load the 'Tidyverse', 2017.
- Zhou, P.-Y. and Chan, K. C.: A model-based multivariate time series clustering algorithm, in: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pp. 805–817, Springer, 2014.