

## A Supplementary Material

### A.1 Dataset Documentation

We provide statistics and plots describing our dataset to give a high level overview of the properties and distribution of the variables. Our tables and plots have been inspired by the format in [37]. In tables 2, 3 and 4 we provide basic statistics on the different nominal, ordinal and continuous variables present in our dataset.

Table 2: Description of Nominal Variables

Name	Type	Count	Unique entries	Most frequent	Least frequent	Missing
deviceId	string	12542183	13	00000000c37f0aa8 (1393279)	000000008f525c6e (427436)	0%
uid	string	12542183	12542183	multiple detected	multiple detected	0%

Table 3: Description of Ordinal Variables

Name	Type	Count	Unique entries	Most frequent	Least frequent	Missing
dateTime	string	12542183	5725406	2021-01-13 20:42:05 (12)	multiple detected	0%

Table 4: Description of Continuous Variables

Name	Type	Count	Min	Median	Max	Mean	Std	Missing	Zeros
lat	number	12542183	28.48655	28.57932	28.72000	28.57855	0.050018	0%	0%
long	number	12542183	77.10014	77.24073	77.32076	77.25250	0.03945	0%	0%
pm1_0	number	12542183	1	110	1730.5	120.34796	57.27231	0%	0%
pm2_5	number	12542183	1	183	1792	207.92479	114.36323	0%	0%
pm10	number	12542183	1	199	1903	226.11056	123.8647	0%	0%

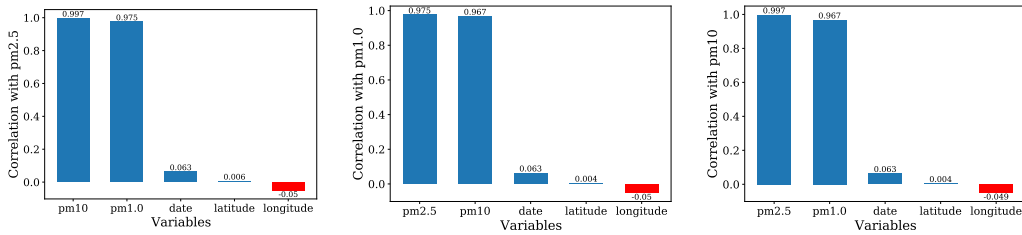


Figure 8: Correlation of different variables with PM2.5, PM1 & PM10 values across the dataset

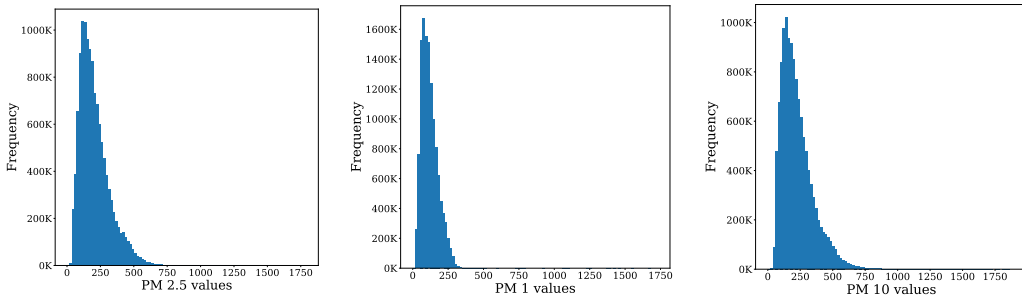


Figure 9: Frequency distribution of PM2.5, PM1 & PM10 values across the entire dataset

Fig 8 shows how the different variables in our dataset correlate with our PM variables while 9 shows the distribution of the PM variables. We see that the three PM sizes of 2.5, 1.0 and 10 are highly correlated in Fig 8, while in absolute values, PM 2.5 and PM 10 see higher magnitudes than PM 1.0 in 9.

Fig 10 (a) shows the distribution of our devices over the entire dataset. It is represented by a bar graph over the different devices. 'f0aa8' is the most frequent device whereas '25c6e' is the least

frequent device. Fig 10 (b)-(d) show our average PM2.5<sup>4</sup> plots for November, December and January separately. We average out all the PM2.5 values for each month, after rounding the latitude and longitudes till the 3rd decimal place, and plot the rounded points in our trajectories. We observe two things here. The first, and most immediate, is that the trajectories differ for all the three months. Since our devices are attached to public buses, the routes often change and this is observed in the plots. The second is that the PM2.5 distribution varies from one month to the next. (b) shows that November has relatively lesser pollution levels across the city with only a small number of regions showing heavy pollution levels. The circles plotted are light in color except around some areas, including the center of our area of study, where the Kushak Nallah bus depot is located. (c) shows December with an increase in the pollution levels with a cluster of really dark points in the center of our area of study. The trajectories change and include a larger area now. Finally, (d) shows that January also has fairly high PM, with it's own regions of densely located dark points. Both December and January indicate the winter season peak and thus, higher average pollution levels are expected.

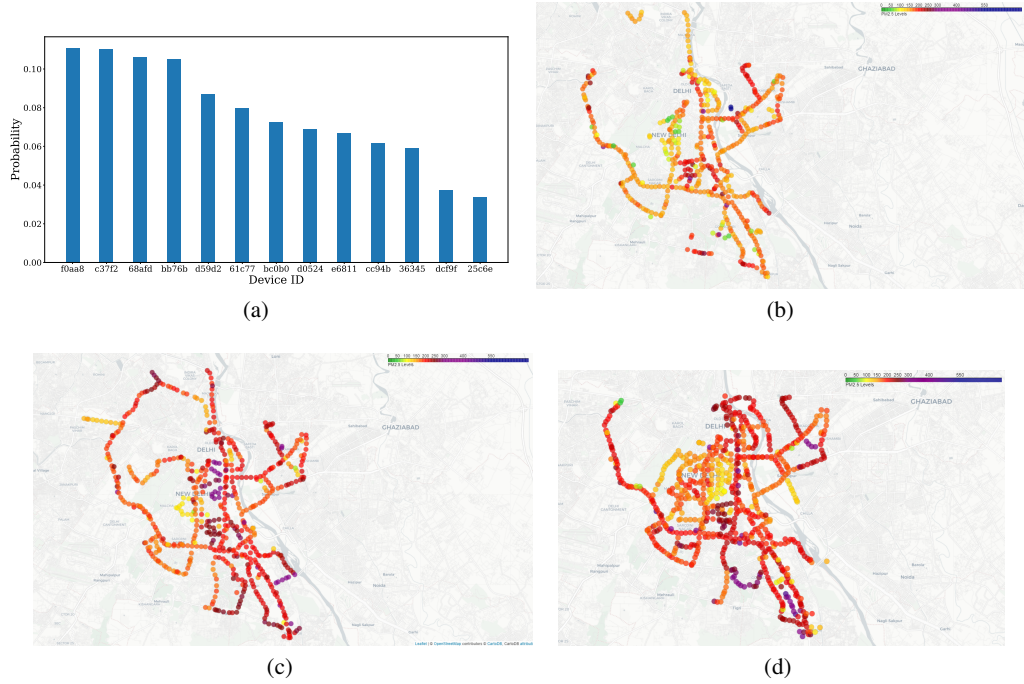


Figure 10: (a) Contribution of individual devices in the entire dataset, in terms of number of records collected. 'f0aa8' is the most frequent device whereas '25c6e' is the least frequent device. (b)-(d) Average PM2.5 maps for each month in our data. These maps show the average PM2.5 values in each month for the points in our trajectories. Lighter colored points indicate lesser pollution levels while darker colored points indicate mode pollution levels. (b) November has relatively lesser pollution levels across the city with only a small number of regions showing heavy pollution levels. (c) December shows an increase in the pollution levels with a cluster of really dark points in the center of our area of study. (d) shows that January still has higher pollution levels.

## A.2 GPS data cleaning details

As discussed in Section 2, our GPS entries are sometimes 0, when our IoT unit cannot lock the satellites through the metallic body of the bus, especially when the bus moves under a bridge or flyover, or passes an underground tunnel. To deal with these, we formulate a problem as follow: suppose we know the location tuples  $(x_1, y_1)$  and  $(x_2, y_2)$  for a bus at timestamps  $t_1$  and  $t_2$ , we wish to find its possible location at time  $t \in (t_1, t_2)$ . We solve this using Map Matching or Linear Interpolation, both of which give good accuracy, with the latter having lower runtimes. We detail the two methods here.

<sup>4</sup>PM values are measured in  $\mu\text{g}/\text{m}^3$

**Map Matching:** Map matching is a commonly performed process to infer a path on a road network from a noisy GPS trajectory [16]. A map-matching algorithm takes recorded serial location points and relates them to edges in an existing street graph (network) in order to represent the trajectory of a user or vehicle. We used Fast Map Matching<sup>5</sup> (fmm), an open source map matching framework in C++ and Python based on [17]. It provides detailed matching information and scales to millions of GPS points and road edges. This helps us achieve two objectives. The first being filling the missing points and the second being correcting any non-missing points in case the recorded values are shifted from the predicted route. After implementing this however, we faced a number of challenges which included the breaking of the algorithm if any two consecutive points were too far, entire trajectories being unmatched if algorithm failed at any point, extensive tuning of the algorithm’s parameters and the algorithm taking a lot of time to compute. Fig. 11 (a) shows the time taken on the data of four different days along with the length of trajectory. For a given graph size, the time taken by map-matching is directly proportional to the number of points in the GPS trajectory to match.

**Linear Interpolation:** Linear interpolation could be easily used in our case because of the high average sampling rate of devices (one ping per 3 sec). For such small distances, linear interpolation serves as a good enough approximation. Also, as linear interpolation consists of simple calculations, it is an extremely fast method (around 1 second; constant in input size) to apply which makes it easy to process the datasets. We fixed a threshold distance,  $\Delta$ , up to which linear interpolation is considered to be valid. To fill any missing point  $P$  recorded at time  $t$ , we did the following:

Identify non-missing points just before and after  $t$ , say  $P_1 : (x_1, y_1, t_1)$  and  $P_2 : (x_2, y_2, t_2)$  respectively. Calculate distance  $d$  between  $P_1$  and  $P_2$ . If  $d$  is greater than  $\Delta$ , leave  $P$  as it is, otherwise, fill  $P$  as  $\left( x_1 + \frac{t-t_1}{t_2-t_1} \cdot d \cdot (x_2 - x_1) \right), \left( y_1 + \frac{t-t_1}{t_2-t_1} \cdot d \cdot (y_2 - y_1) \right)$

**Evaluation:** For our data, with  $\Delta$  as 200m, more than 95% of the missing points, on an average, could be filled. We designed an experiment to compare map-matching and linear interpolation. From the list of non-missing points, 10-30% GPS points were randomly sampled and kept as test points. Using the remaining 70-90% known points, these test points were predicted using the two methods and the root mean square error (RMSE) and the mean absolute error (MAE) [33] of the Haversine distance between the predicted coordinates and the true coordinates were computed. The results of the experiment are shown in Fig. 11 (b). It can be seen that linear interpolation performs better than fmm. Hence, we decide to move ahead and fill the incorrect points with linear interpolation.

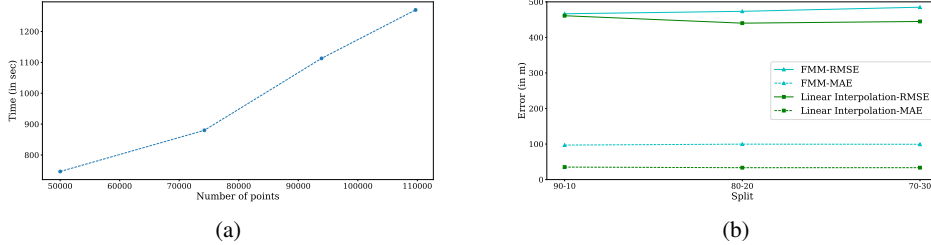


Figure 11: Results of the two interpolation mechanisms. (a) shows the time taken by map-matching on four different days. Map-matching takes a lot of time and thus isn’t useful for where time is a constraint. (b) shows the distribution of errors (RMSE and MAE) for different filled-missing points splits for both the methods. (Threshold  $\Delta = 300$ m for linear interpolation)

### A.3 Spatio-temporal Interpolation Baselines

The different spatio-temporal interpolation baseline models have been described below. The description includes the preprocessing, parameter and hyperparameter settings.

**Preprocessing:** The preprocessing was kept model specific where it could include normalization of the features and target variables, downsizing of the train dataset and finally, creating a graph using the dataset. These were experimented upon using the datasets from different arbitrary days and the best settings were kept for the benchmarking process. The preprocessing details for each model have been mentioned as a part of the model details below.

<sup>5</sup><https://github.com/cyang-kth/fmm>

**GPR:** The Gaussian Process (GP) regression model is defined by specifying the mean and the covariance function the GP prior [28]. While the mean function largely models the trend of the mean values in the model, the covariance function decides how changes in the features affect the values of the target variable. The covariance function is further specified using a kernel function that learns the covariance between the different target values based on the feature values. The covariance or kernel function thus captures the relationship between the PM values and how they vary with changes in our spatial and temporal features. The most common kernels used in a GPR model are RBF, Matern-3/2 and Matern-5/2 [28]. The kernel function has some free hyperparameters that are learnt while optimising the GPR model. The standard hyperparameters are the variance and lengthscale of the kernel along with the variance of the noise in the model.

Firstly, a downsized version of the train dataset is created to facilitate tractable training of the model. The downsized train dataset is created by taking the mean of the PM2.5 variable for all observations which are at a spatial distance of approximately 150m to each other and at a temporal distance of 15 minutes, making this averaged point one single point. This downsized dataset is further normalized, in order to make the mean zero and standard deviation 1 for the input and output features. While the mean function was specified as a constant function, the kernel function was the product of a Periodic and Matern-1/2 kernel. This was done in order to give more importance to the points which were both spatially and temporally close while also allowing us to model the datetime feature as a sine wave. We used ARD, automatic relevance determination (ARD) [28] to allow different values of hyperparameters to be fitted for each feature. For both our Gaussian process-based models we have used [38] to code them.

**Variational GPR:** The Variational GPR allows us to use inducing points to learn a more easily tractable final posterior distribution, by computing the variational lower bound (or ELBO) rather than the GP marginal log likelihood, as compared to the above mentioned GPR model. This allows us to extend the GPR framework to big data. The preprocessing for this model is similar to the preprocessing for GPR. The mean function is kept as a constant function again though the kernel function here is just a Matern-1/2 function with ARD. Further, the inducing points here are limited to 2000 points which help us strike a good balance between computational requirements and accuracy of the model.

**GraphSAGE:** This model has been inspired by GraphSAGE [30]. We aim at learning universal weights, similar to GraphSAGE, which will signify the importance of a neighbour based on some known node values and edge weights. Here we define node values as the value of the pollutant PM2.5 while the edges are created using latitude, longitude and datetime features. It is important to note that our model doesn't use feature values for a point as part of our node definition. Rather, only the pollutant value is defined as the node value and the correlation in pollutant values based on our changes in feature values helps us in predicting our final node value for unknown nodes. Firstly, a graph is created from a downsized version of the train dataset. The downsized train dataset is created by taking the mean of the PM2.5 variable for all observations which are at a spatial distance of approximately 150m to each other and at a temporal distance of 15 minutes and making this one point represent one single node.

An edge is then created between two different nodes only if both the nodes lie within 2 hours of each other temporally and around 2km away from each other spatially. The weight function for aggregation of edge was a polynomial encoding an approximate form of dependence of distance and time between any two nodes with the power of the terms being considered as hyperparameters. Weights of the edges were inversely dependent on the product of a term consisting of the squaring of the haversine distance between those nodes and the temporal distance between them (which was taken as 15 minutes here). During edge formations, to resolve the issue of node isolation in cases of extremely distant nodes, we set the minimum edges to 4 so as to prevent edge scarcity.

Once formed, the graph then goes through two graph-based layers to learn the required weights where embeddings are learnt using the max and mean aggregation layers. These are similar to how they are defined in the original GraphSAGE model. These include the mean aggregator and max aggregator too, as seen in the original GraphSAGE model. After the two layers, there are 3 fully connected neural network layers to predict the final pollutant value. The parameters of these layers are then learnt when this model is trained using the graph created earlier. Finally, instead of just using one graph to train our model, we construct 50 such smaller subgraphs by dividing the main graph randomly into smaller subsets. This makes our training and validation process more robust

as we perform a random walk which protects us against bias in our model. The loss function is the common Mean Squared Error loss between the predicted value of pollutant of all the nodes. For both this model and the Meaner, we have used [39] to code them.

**Meaner:** Meaner simply uses the graph described in our GraphSAGE model and averages the values of the known neighbours of any given unknown node to predict its value. Thus, a weighted mean based on the value of the edges and the nodes neighbouring the unknown node is performed where the weights are the values of the edges.

**ANN:** Extensive testing was performed with various multilayer perceptron models with varied hyperparameters to find the best architecture. We had three input features corresponding to latitude, longitude and datetime. Target variable was PM2.5 value. All the hyperparameters were tuned based on the train and validation datasets.

For preprocessing, we found the mean of PM2.5 values for all observations which were temporally 15 minutes apart. Then we performed normalization of input and output features in order to make their mean zero and standard deviation 1. For the model, a 6-layered architecture was found to be the best (one input layer, one output layer and 4 hidden layers). All the hidden layers had 200 densely connected hidden units each and weight initialization of He et al. [31] was used. Further, Rectified linear unit (ReLU) was used as the activation function for all layers. For optimization, we used Adam [32] optimizer with an exponentially decaying learning rate. Finally, the loss function used was Mean Squared Error (MSE) loss between predicted and actual PM2.5 values. Model implementation was done using open source Pytorch library [40].

#### A.4 Anomaly metrics computation details

As discussed in Section 5, we compute six anomaly metrics on our dataset. Each of these metrics needs several thresholds, which we find by running some heuristics to compute the relevant statistics from the dataset. The heuristic to compute the first anomaly metric of **samples recorded per minute** is given in Algorithm 1. This metric checks for faulty devices which might be sampling more or less than expected rate. Fig. 5a in the paper shows ideal samples collected per minute should be around 20. If it deviates too much, that device is anomalous. The amount of deviation allowed is calculated statistically by observing the distributions for several days. Our algorithm finds the upper bounds and lower bounds of the median ( $\Theta_{50}^L, \Theta_{50}^U$ ), 25<sup>th</sup> percentile ( $\Theta_{25}^L, \Theta_{25}^U$ ) and 75<sup>th</sup> percentile ( $\Theta_{75}^L, \Theta_{75}^U$ ) of the expected distribution.

---

##### Algorithm 1 Finding Metric 1 parameters

---

**Step 1 :** Define a set of days(n) based on which we want to determine the thresholds (X). These days can be manually picked by observing the box plots as in figure 5a.

**Step 2:** For each day and each device we collect the number of readings from active minutes.

**Step 3:** From the set of distributions(n), we collect median, 25<sup>th</sup> percentile and 75<sup>th</sup> percentile of each distribution. Lets denote it by  $D_{50}, D_{25}, D_{75}$  respectively.

**Step 4:** To find a considerable range of these parameters we choose 50 percent of central data ruling out first and last 25 percent of data as outliers.

<b>Step 5:</b> $\Theta_{50}^L = 25^{th}$ percentile of $D_{50}$	$\Theta_{50}^U = 75^{th}$ percentile of $D_{50}$
$\Theta_{25}^L = 25^{th}$ percentile of $D_{25}$	$\Theta_{25}^U = 75^{th}$ percentile of $D_{25}$
$\Theta_{75}^L = 25^{th}$ percentile of $D_{75}$	$\Theta_{75}^U = 75^{th}$ percentile of $D_{75}$

---

Other metrics use similar heuristics. Fig. 12 shows the summary of anomalies detected based on these heuristics and the thresholds listed in Table 5. Fig. 12(a) shows the summary for metrics 1 and 2, which compute **samples recorded per minute** and **number of minutes active in an hour** respectively. Both these metrics show anomalies in the early phase of the mobile sensor deployment, and the anomalies become rarer after Dec 10<sup>th</sup>, 2020. This is similar to the flat part parallel to x-axis in Fig. 6(b) in the paper, which showed **number active hours in a day** i.e. metric 3. The early phase of IoT deployment saw issues in 4G radio signals based on the placement and packaging of the unit, which was changed until data collection became smooth and anomalies for these metrics became rare. The remaining few anomalous points for the rest of Dec, 2020 and Jan 2021, resulted from occasional 4G signal losses, which is common in urban tunnels in Delhi as the bus travels. 1-2 units consistently performed poorly, which have been debugged to have power supply issues from the bus.

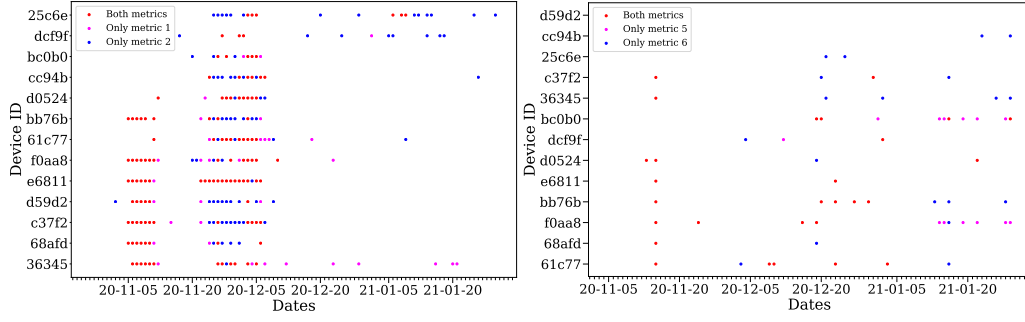


Figure 12: Anomalies detected by our heuristics on the whole dataset (Nov 1<sup>st</sup>, 2020 to Jan 31<sup>st</sup>, 2021). Presence of marker represents anomalous device for corresponding date. (a) shows combined anomalies for metrics 1 and 2. (b) shows combined anomalies for metrics for metrics 5 and 6.

Fig. 12(b) shows the summary anomalies for metrics 5 and 6, which compute **inter-sensor** and **intra-sensor** differences in recorded PM values. We again used the detected anomalies to check the devices and the deployment partner (bus company). In most cases, a deviation of a sensor from other sensors happened when there was local electrical maintenance work like soldering ongoing in the bus, a common phenomena when the buses are parked in the depot at night. 1-2 units had some persistent problem towards the end of month Jan, which on checking were found to accumulate some dust near the fan that suck in air. Some cleaning of the units is therefore recommended, which can be done during regular maintenance of the buses in the depot.

Based on our experience of running a live IoT network on public buses in a developing country, such anomaly detection is vital for quick debugging and fixing of issues. These plotted anomalies are therefore included as csv files in our dataset, so that ML researchers interested in automated anomaly detection methods for IoT networks, can use these as reference ground truth. The heuristic codes are also open-sourced, so that the thresholds can be adjusted based on the ML researchers requirements of more or less conservative anomaly detection.

The summary anomalies are represented in one-hot encoded format in csv files. In all the files (except for metric 4), the first column contains the dates for three months and the first row contains all the 13 device Ids. Presence of 1 indicates anomaly for corresponding date and device. In the file of metric 4, the first row represents the 16 regions in consideration rather than the device Ids.

<sup>6</sup>Please refer Algorithm 1 for detailed meanings of  $D_{25}$ ,  $D_{50}$ ,  $D_{75}$ . The parameter values are found based on date Jan 16 to Jan 31 2021

<sup>7</sup>Please refer section 5, Anomaly metric 4 for detailed description of the symbol and calculation of its value.

<sup>8</sup>Please refer section 5, Anomaly metric 5 & 6 for detailed description of the symbols and calculation of their values.

Table 5: Thresholds used in computation of the six anomaly metrics

Symbol	Metric	Value	Description
$\Theta_{25}^L$	1	19	Lower bound of sampling rate, 25 <sup>th</sup> percentile of $D_{25}$ <sup>6</sup>
	2	20.0	
$\Theta_{25}^U$	1	22	Upper bound of sampling rate, 75 <sup>th</sup> percentile of $D_{25}$ <sup>6</sup>
	2	38.5	
$\Theta_{50}^L$	1	20	Lower bound of sampling rate, 25 <sup>th</sup> percentile of $D_{50}$ <sup>6</sup>
	2	38.0	
$\Theta_{50}^U$	1	23	Upper bound of sampling rate, 75 <sup>th</sup> percentile of $D_{50}$ <sup>6</sup>
	2	59.0	
$\Theta_{75}^L$	1	21	Lower bound of sampling rate, 25 <sup>th</sup> percentile of $D_{75}$ <sup>6</sup>
	2	57.6	
$\Theta_{75}^U$	1	24	Upper bound of sampling rate, 75 <sup>th</sup> percentile of $D_{75}$ <sup>6</sup>
	2	60.0	
$\gamma$	3	1(static sensor) 800(mobile sensor)	Minimum data points in an hour to be termed as Active hour
$\tau$	3	10	Minimum number of active hours for which a sensor must report data to be termed as ideal for the day
$\Theta_{15}$	3	variable	15 <sup>th</sup> percentile of distribution of active hours for a sensor in the previous 15 days
$\delta$	4	41	Maximum allowed percentage difference of a region with its moving average <sup>7</sup>
$\Theta_{25}$	5	variable	25 <sup>th</sup> percentile of distribution of PM values of a device during an hour
	6		
$\Theta_{75}$	5	variable	75 <sup>th</sup> percentile of distribution of PM values of a device during an hour
	6		
max_IQR	5 6	39.77	Maximum allowed IQR, 90 <sup>th</sup> percentile of distribution of all IQRs in training set <sup>8</sup>
buffer	5	105.45	Used for calculating the majority PM range <sup>8</sup>
	6	46.36	