

A Scalable Deep Learning System for Monitoring and Forecasting Pollutant Concentration Levels on UK Highways

Abstract

The construction of intercity highways by the government has resulted in a progressive increase in vehicle emissions and pollution from noise, dust, and vibrations despite its recognition of the air pollution menace. Efforts that have targeted roadside pollution still do not accurately monitor deadly pollutants such as nitrogen oxides and particulate matter. Reports on regional highways across the country are based on a limited number of fixed monitoring stations that are sometimes located far from the highway. These periodic and coarse-grained measurements cause inefficient highway air quality reporting, leading to inaccurate air quality forecasts. This paper, therefore, proposes and validates a scalable deep learning framework for efficiently capturing fine-grained highway data and forecasting future concentration levels. Highways in four different UK regions - Newport, Lewisham, Southwark, and Chepstow were used as case studies to develop a REVIS system and validate the proposed framework. REVIS examined the framework's ability to capture granular pollution data, scale up its storage facility to rapid data growth and translate high-level user queries to structured query language (SQL) required for exploratory data analysis. Finally, the framework's suitability for predictive analytics was tested using fastai's library for tabular data, and automated hyperparameter tuning was implemented using bayesian optimisation. The results of our experiments demonstrated the suitability of the proposed framework in building end-to-end systems for extensive monitoring and forecasting of pollutant concentration levels on highways. The study serves as a background for future related research looking to improve the overall performance of roadside and highway air quality forecasting models.

Keywords: Urban Air Pollution, Air Quality Prediction, Highway, Deep Learning, Big Data, Internet of Things

1. Introduction

Long-term exposure to air pollution is the most significant environmental threat to human health (Public Health England 2019). According to World Bank (2022), the global cost of the adverse health effects associated with exposure to air pollution is \$8.1 trillion, equivalent to 6.1 per cent of global GDP. It is, therefore, surprising that a substantial fraction of the UK populace (particularly those that commute to their various destinations via highways) are still susceptible to the adverse health effects of air pollutants along the UK highways (Vohra et al. 2021). Due to exposure to motor vehicle exhaust emissions, non-exhaust related pollution from brake and tyre wear, and particles from highway construction (Barikayeva et al. 2018), commuters are constantly at risk of high concentrations of air pollutants (e.g., $PM_{2.5}$, PM_{10} , NO_2). These pollutants are some of the most life-threatening road pollutants, which

12 have been linked to cardiovascular and respiratory illnesses (Mabahwi et al. 2014, Alvanchi
13 et al. 2020). According to Public Health England (2019), between 2017 and 2025, these air
14 pollutants will cost the NHS and social care system in England a total of £1.6 billion.

15 Hence, there is a pressing and cogent need to find innovative and sustainable ways to mon-
16 itor air pollutants and curb their devastating effects on health and human capital, as well
17 as associated GDP losses (DEFRA 2020). According to (Alvanchi et al. 2020), monitoring
18 particulate matter ($PM_{2.5}$, PM_{10}) and other highway pollutants like NO_2 is not a straight-
19 forward task because pollutants tend to decay and diffuse into the background concentration
20 within 200m from the source. Furthermore, highway speed limits and traffic congestion com-
21 plicate things further as they result in varying driving patterns such as sudden slow-downs
22 and speedups, which elevate these pollution levels or limit their dispersion (Karner et al.
23 2010, Zhang & Batterman 2013). In response to the ever-increasing impacts of air pollution
24 and its associated intricacies, the UK government has invested about £100 million to proac-
25 tively tackle air quality (AQ) challenges to protect health and support clean growth (DEFRA
26 2019). However, despite these investments by the UK, the issue of how to proactively tackle
27 and ultimately improve air quality across UK highways persists.

28 According to Barthwal & Acharya (2018), most countries monitor air pollution using sta-
29 tionary monitoring stations operated by government authorities. Figure 1 illustrates how the
30 UK currently monitors highways to come up with its ultra low emission policies. Highways
31 are monitored by Highways England (a government-owned company charged with operat-
32 ing, maintaining, and improving motorways in England) via its automatic urban and rural
33 network (AURN), which collects sparse air pollutant data. However, evidence suggests that
34 these air quality analysers are relatively heavy and expensive to install or maintain (Carullo
35 et al. 2007, Barthwal & Acharya 2018). Therefore, it is impracticable for Highways Eng-
36 lands’ monitoring stations to be deployed across the UK to capture pollutant concentration
37 levels and improve air quality. On the other hand, low-cost/off-the-shelf IoT sensors that
38 have been proposed in previous studies (e.g., Badura et al. (2018), Borghi et al. (2018),
39 Budde et al. (2018)) for monitoring air quality are plagued with interference issues from
40 weather, cross-sensitivities between pollutants and ageing effects of integrated circuit tech-
41 nology (Karagulian et al. 2019). These limitations are coupled with the fact that it would
42 take a significant and environmentally unfriendly investment to install low-cost IoT sensors
43 across every road in the UK.

44 Asides from inefficient highway air quality monitoring, another major challenge rests on
45 the issue of how data disparity and isolated data sets affect the accurate prediction of pollu-
46 tant concentration levels. Evidence suggests that data sources that could collectively predict
47 pollutant concentration levels (e.g., historic pollution, GIS location data, traffic flow, weather
48 and background pollution) exist in silos, thus resulting in numerous integration and data pro-
49 cessing issues (Umadevi & Geraldine Bessie Amali 2020). In addition, existing forecasting
50 methods are riddled with computational issues like scalability and memory demand which
51 limits their optimal adoption (Zhang et al. 2012). For instance, Alléon et al. (2020) and
52 Lee et al. (2020) attempted to develop large-scale air quality forecasting systems. However,
53 the authors highlighted the inability to integrate additional granular data and insufficient

54 computational power as the shortcomings of their research.

55 On the back of significant advancements in scalable machine learning (ML) approaches
56 such as deep learning, which are known to thrive on huge data (Akinosho et al. 2020), this
57 study, therefore, proposes a scalable deep learning framework for monitoring and forecast-
58 ing pollutant concentration levels on UK highways. This framework leverages internet of
59 things (IoT) sensors for real-time monitoring, graphics processing units (GPUs) for parallel
60 computing, big data for scalable storage and deep learning for forecasting highway pollu-
61 tant concentration. In the design of a system that implements the proposed framework, the
62 following objectives were set for this study:

- 63 • Develop, calibrate, and deploy energy efficient hardware devices with practicable accu-
64 racy to capture real-time pollution data on four different UK highways.
- 65 • Integrate missing or inaccurate data from heterogeneous sources to enhance forecasting
66 accuracy of the developed model.
- 67 • Develop and evaluate a baseline deep learning model to make hourly predictions of
68 PM2.5, PM10 and NO2 concentration levels due to the deadly nature of these pollu-
69 tants.
- 70 • Perform a system scalability test to determine response time and throughput as hard-
71 ware device load increases.

72 This manuscript is structured as follows: Immediately after introduction follows a sec-
73 tion that summarises the research methodology adopted to achieve the highlighted research
74 objectives while section 3 highlights features of the proposed framework. Afterwards, the
75 development process of a prototype system that implements the proposed framework is dis-
76 cussed in section 4. The discussion also includes the scalability test that was performed on
77 the system. Section 5 correlates the findings of this study with existing research and high-
78 lights its relevance to practice. Finally, conclusions are drawn, and future research directions
79 are indicated in Section 6.

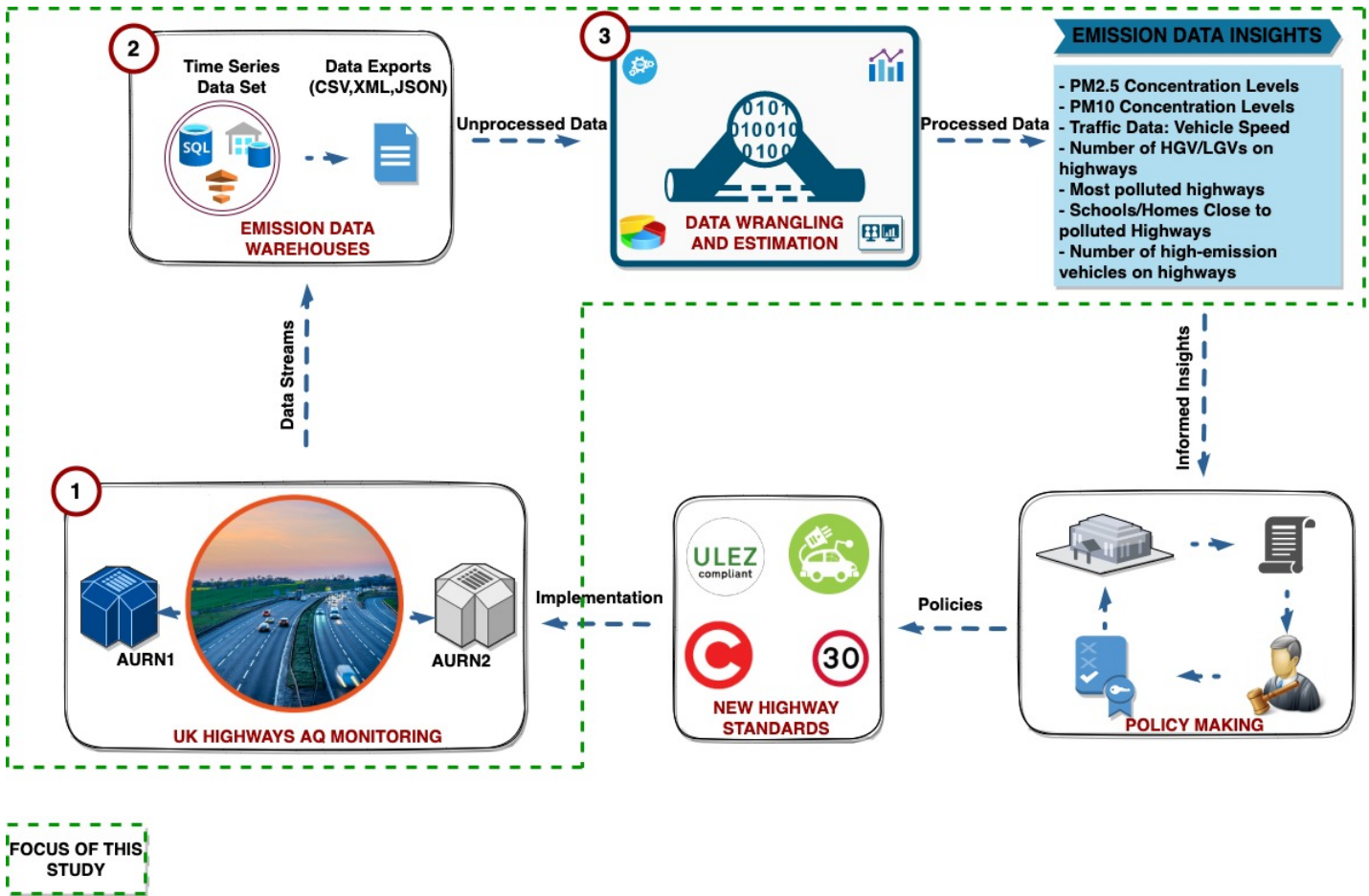


Figure 1: The current situation of highway AQ monitoring in the UK. This study seeks to address three main challenges which include: 1) expensive cost of deploying monitoring stations such as the AURN on highways 2) Data silos/segreated data operated by different data agencies 3) Inefficient air quality estimation methods

80 2. Research Methodology

81 A combination of experimental design with case-study methodology was adopted to
82 achieve the identified research objectives of this study. The rationale behind this approach
83 is based on the need to fulfil two principal goals. One, system implementation is needed to
84 demonstrate the practicality of the proposed framework using a developed system. This ap-
85 proach reflects an experimental design method of research. On the other hand, the adoption
86 of case-study strategy is to test the results of the developed system in disparate real-life envi-
87 ronments. This approach is quite effective and has been utilised in related research involving
88 multiple case-studies (Chen et al. 2015, Zheng et al. 2015). The experimental design approach
89 was also used to conceive the layered architecture of the proposed framework with each layer
90 addressing at least one research objective. Layering remains a prevalent application design
91 technique that allows the disintegration of a complex software system into modules. Lay-
92 ers within the proposed framework consist of libraries, programming languages, and services
93 required for monitoring and forecasting.

94 A careful market analysis of “grey literature” and an extensive review of academic pub-
95 lications revealed several sensors suitable for designing monitoring units to address the first
96 research objective. Google scholar and scientific databases such as Scopus and ScienceDirect
97 were also used to search for academic publications, while Google’s search engine revealed
98 additional sensor manufacturers. We limited our search to the three pollutants of interest
99 in this study - NO_2 , $PM_{2.5}$ and PM_{10} . A similar search of relevant integration libraries
100 and big data frameworks informed the framework’s approach to solving integration and data
101 storage challenges. From the array of available options, enterprise frameworks that allow the
102 integration of data from legacy as well as newly built systems were selected. There are a
103 lot of algorithms that are available for air quality forecasting. However, it was important to
104 choose a scalable machine learning approach such as deep learning that has shown significant
105 promise using distributed computing clusters (Sergeev & Del Balso 2018, Chen et al. 2019).

106 **3. A Proposed Deep Learning Framework for Highway AQ Monitoring and Pre-** 107 **prediction**

108 The proposed framework is a four-layered architecture composed of the hardware layer,
109 data storage layer, integration layer and analytics layer as depicted in Figure 2. This section
110 introduces these layers and their functionalities.

111 **3.1. Hardware Layer**

112 This layer serves as the entry point for the entire framework. It initiates the monitoring
113 and analytics process by ensuring that real-time data are captured and subsequently trans-
114 ferred to a cloud platform for data aggregation. A typical real-time sensing device in this
115 layer would push data at an interval of 30secs-1min and be able to sense multiple pollutants
116 and capture weather data. Other device functionalities such as self-powering capability, edge
117 computing and on-board intelligence are desirable but not entirely mandatory for monitoring.
118 Multiple gateways and a cloud platform are essential for this layer to function as required.
119 The cloud platform will store captured data, but on-device storage will also be helpful to
120 avoid data loss when data transfer fails. Additional data on vehicle categories and traffic
121 flow in this layer will provide more insights into the 'culprit' vehicle that contributes the
122 most to highway pollution. Advanced computer vision and edge computing technologies can
123 enable this functionality in monitoring devices through embedded ML models. Development
124 technologies relevant to this layer include VHDL, Verilog, FPGA, and Arduino.

125 **3.2. Data Storage Layer**

126 This layer stores pollution data and model weights. Readings captured from deployed
127 sensing devices are either sent immediately to this layer or stored temporarily and pushed
128 later through HTTP post requests. The data storage layer is responsible for ensuring data
129 consistency, security and integrity. According to Ahmed et al. (2017), it is best practice to
130 have the unified prediction service (UPS) reside close to the historic pollution data to reduce
131 latency. Hence, this layer also houses weights and parameters from training pollutant con-
132 centration forecasting models. Data stored in this layer are bound to increase exponentially,
133 and necessary technologies to configure big data storage must be put in place. Relevant tech-
134 nologies such as hadoop, spark and hive are possible open-source options to consider in this
135 configuration. Data streaming frameworks like Apache Kafka or ActiveMQ are also available
136 for real-time sensing of changes in this layer and to send alerts in the event of data trans-
137 fer failures. Triggers, procedures and packages are useful to automate most of the required
138 database tasks such as populating tables, generating logs or automatically generating SQL
139 for data aggregation.

140 **3.3. Integration Layer**

141 The data integration layer ingests data from third-party sources into a central repository.
142 The layer handles this data ingestion using the extract, transform and load (ETL) process.
143 External data can include pollution data captured by other monitoring stations, highway
144 geographical data, meteorological data and traffic data. The essence of this layer is to ensure
145 that data not captured in the hardware layer by the monitoring devices can be integrated
146 into the system to improve the performance of developed estimation models. If the suggested

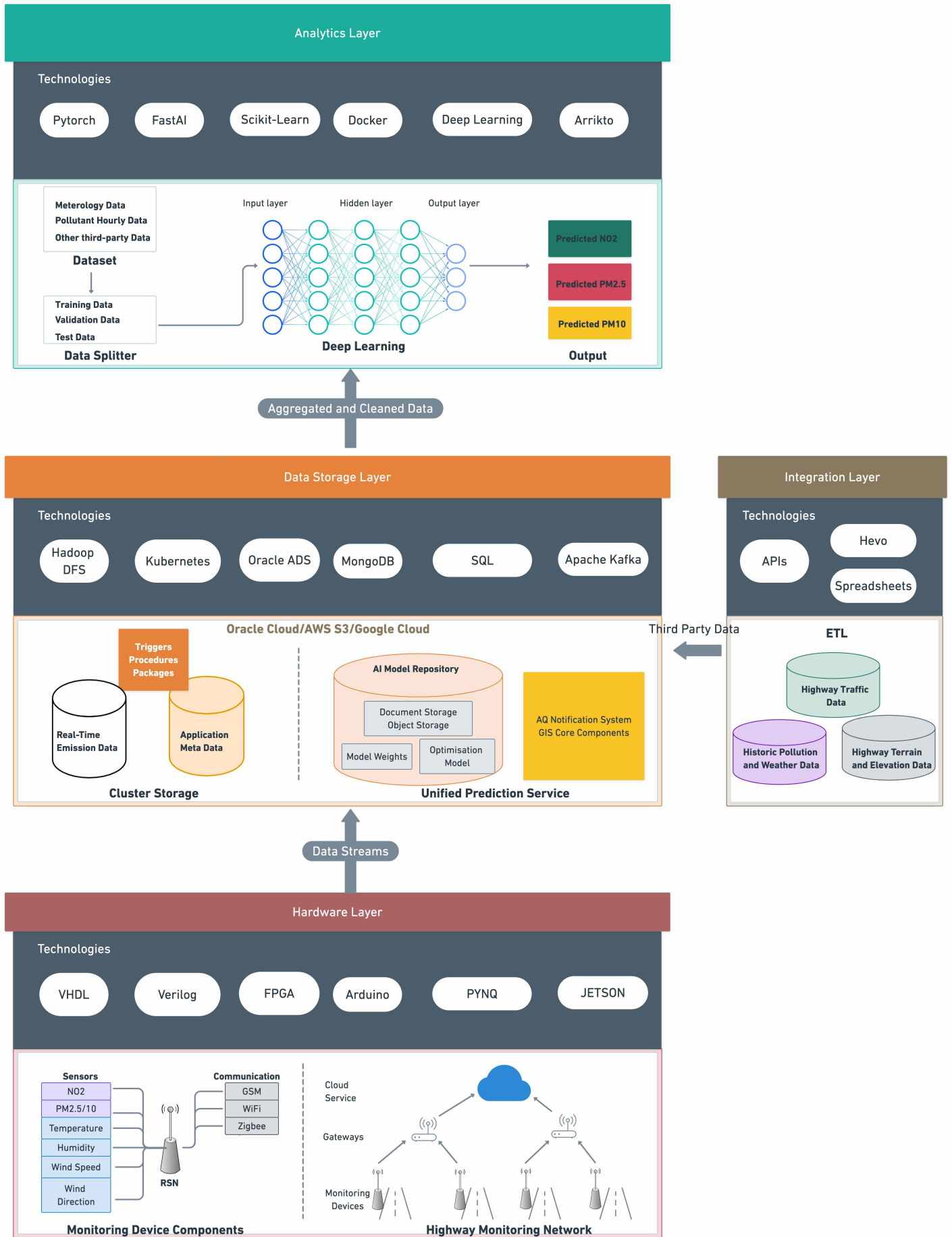


Figure 2: Deep Learning Framework for Highway Air Quality Monitoring and Prediction

147 functionalities of the hardware layer are too expensive to implement, this layer can grab
148 open-source or paid data from available online sources. Data can be downloaded in different
149 formats such as TXT, JSON, XML and CSV or exposed as external links. The data from this
150 layer should be stored as separate tables in the data storage layer for unique identification
151 and also to avoid mix-ups with existing data.

152 3.4. Analytics Layer

153 The analytics layer handles exploratory and inferential analysis of historic highway pol-
154 lution data to estimate future air quality. The layer extracts data from the data storage
155 layer for model training and validation. Essential data pre-processing steps such as data con-
156 sistency verification, target attribute transformation, feature extraction, data encoding and
157 data imputation are carried out in this layer as part of the first stages of training. A machine
158 learning approach suitable for tabular or time-series data such as the historic pollution data is
159 required for estimation. Deep learning is one of many machine learning approaches that has
160 stood the test of time (Akinosho et al. 2020). Frameworks and libraries such as fastai, scikit-
161 learn, PyTorch and TensorFlow make it relatively easy to train a baseline model. Additional
162 functionalities that are beginning to gain traction and could be included in implementing this
163 layer is MLOps - model maintenance in the production environment. MLOps encompasses
164 automation and monitoring steps such as continuous integration, deployment and training
165 on data collected in production.

166 4. Development and Deployment of the REVIS System Prototype

167 In this section, the proposed framework is validated for practicality through the imple-
168 mentation of a Real-Time Highways Emission Visualisation (REVIS) platform use case. The
169 framework was tested for scalability and performance through different stages of data collec-
170 tion, exploratory data analysis and predictive model development.

171 4.1. REVIS Highway Monitoring Devices

172 The development and evaluation steps of the monitoring devices and the deployment
173 strategy adopted are highlighted in this section.

174 4.1.1. REVIS Device Development and Evaluation

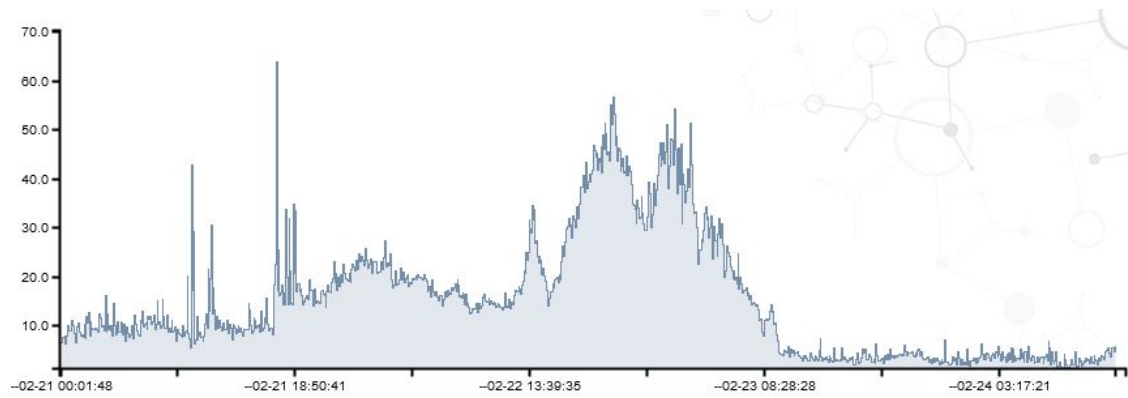
175 REVIS demonstrates the hardware layer through the development and calibration of de-
176 vices with built-in sensors to measure the atmospheric composition of NO_2 , $PM_{2.5}$ and PM_{10} ,
177 alongside weather parameters - pressure, temperature and relative humidity. Table 1 below
178 summarises details of manufacturers of the chosen sensors and their accuracy figures. Each
179 REVIS device required an excellent design of both analogue and digital circuitry around it
180 and several stages of calibration. The Alphasense NO_2 sensor for example, showed during
181 experimentation that it was best suited for fixed sensing installations and urban air monitor-
182 ing since varying meteorological conditions had a significant influence on it's readings. The
183 sensor's cross-interference with the $PM_{2.5}$ SPS30 sensor and detection range limits (DRL)
184 were also evaluated using equation 1

$$DRL = 3.3\sigma/S \tag{1}$$

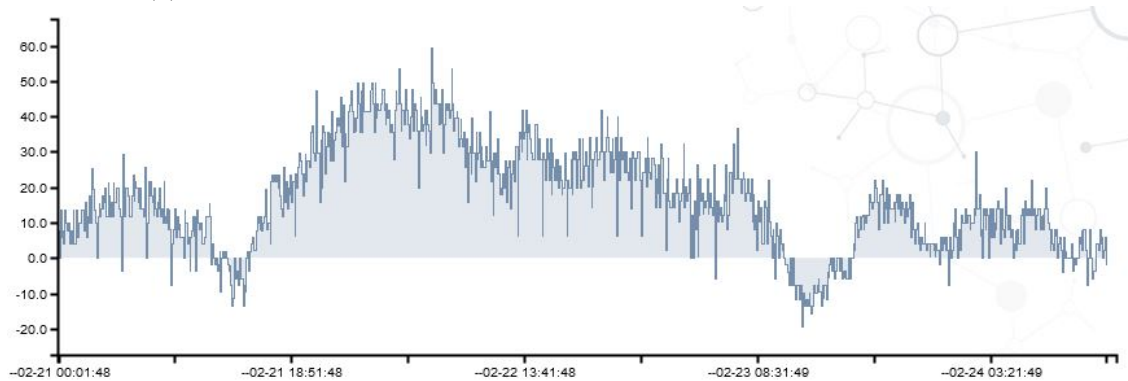
185 where S denotes the calibration curve's slope, and σ denotes the standard deviation of the
 186 sensor response in the absence of air (Shrivastava et al. 2011). The nearest AURN stations
 187 to the monitoring devices were identified for field evaluation. The selected stations were
 188 deemed suitable for calibration since they were close to deployed sensors and mainly provided
 189 missing weather data and also average hourly measurement of the pollutants of interest. Data
 190 from the REVIS devices were averaged over an hour for appropriate comparison with the
 191 reference data. Figure 3 shows $PM_{2.5}$ and NO_2 readings on one of the REVIS devices after
 192 calibration. Aside from the occasional underestimated measurement of the NO_2 sensors,
 193 other sensors such as $PM_{2.5}$ and PM_{10} showed close estimates to the reference measurements
 194 with correlation coefficient $r > 0.8$.

Table 1: Sensor Specifications and Accuracy

Measured Quantity	Units	Sensor used	Accuracy	Comments
Temperature	$^{\circ}C$	Texas: HDC2010	± 40	Could be affected by direct sunlight, depending on how well airflow works within the unit - may require additional physical shading.
Relative Humidity	%	Texas: HDC2010	± 3 start of life $\pm 0.25/\text{yr}$ drift	As above
Pressure	hPa	ST: LPS22HB	± 1	
$PM_{2.5}$ and PM_{10}	$\mu g/m^3$	Sensirion: SPS30	$\pm 10 \mu g/m^3$ $\pm 10\%$	Over 0-100 $\mu g/m^3$ range Over 100-1000 $\mu g/m^3$ range
NO_2	ppb	Alphasense: B43F	NO2- Approx. ± 20	Careful design and several stages of calibration are required when measuring tiny gas concentrations



(a) $PM_{2.5}$ field readings after calibration



(b) NO_2 field readings after calibration

Figure 3: Calibrated NO_2 and $PM_{2.5}$ readings from field. Vertical units are in $\mu g/m^3$ for $PM_{2.5}$ and ppb for NO_2 . Even with the calibration, NO_2 readings sometimes record negative readings because of temperature and humidity effects.

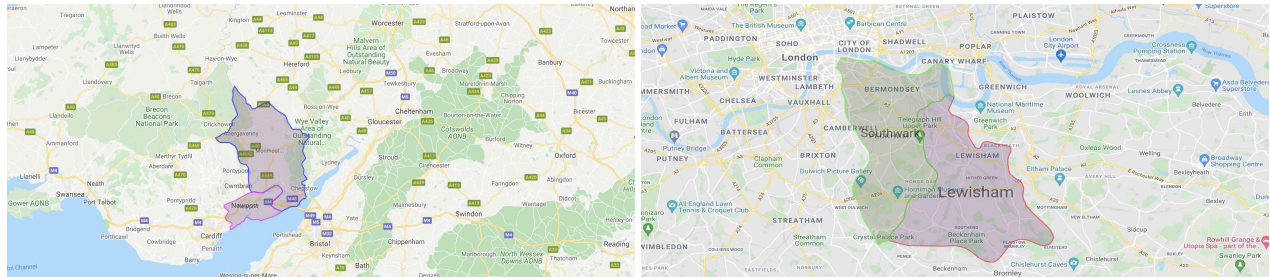
195 4.1.2. Device Deployment in Case Study Regions

196 Major highways in London, Newport and Chepstow were chosen as case studies for this
197 research. The city of London is made up of 9 million inhabitants which includes 4.49 million
198 males and 4.51 million males (ONS 2020). With 74.9% of this population belonging to the
199 working age 16-64 years (ONS 2021), the government faces the challenge of addressing traffic
200 congestion hurdles across the city. A survey reported in TFL (2019) showed that 59% of
201 Londoners tend to use the bus at least once a week, with car passenger commuting more
202 popular among the younger populace. Newport and Chepstow are located in the south
203 eastern region of Wales. According to ONS (2018), 1.53 million residents live in the region
204 with a population density of 546 person/ km_2 . The region also has the highest and lowest life
205 expectancy figures across wales with Monmouthshire boasting the highest life expectancy.
206 74.3% of the employed populace prefer to journey by motorcycle, van or car while 8.8% choose
207 to travel by bus or train (Statswales 2020). Major highways in the region such as the M4,
208 A48 and A466 highways connect neighbouring cities.

209 Figure 4 depicts the distribution of REVIS devices in these cities. In London, twelve
210 devices were distributed on sections of the A302, A2209 and A1203 highways. One device
211 was placed 92.79m from Junction 25 of the M4 highway in Newport and another device
212 was positioned close to The A48 motorway in Chepstow. The deployment approach that was
213 adopted during the distribution of these devices ensured three key requirements: (1) sufficient
214 highway length (2) cellular data connectivity and (3) electrical/solar power availability. It
215 was also necessary that device installation required minimum technical skills and data was
216 captured for a minimum of 6-8 months.

217 4.2. Exploratory Analysis of Pollution and Weather Data

218 It is important to verify data consistency before commencing model training in ML re-
219 gression tasks such as the one being considered. The minimum recommendation is to confirm
220 the total number of rows and columns within the data, as this may have been compromised
221 during data transfer (Bilal & Oyedele 2020). This section analyses the impact of weather pa-
222 rameters and the case-study region on pollutant levels. Although data was captured between
223 November 2020 and August 2021, missing data in the early stages of deployment (shown in
224 Figure 5 below) influenced the decision to analyse data between February 2021 and August
225 2021 when missing data was minimal.

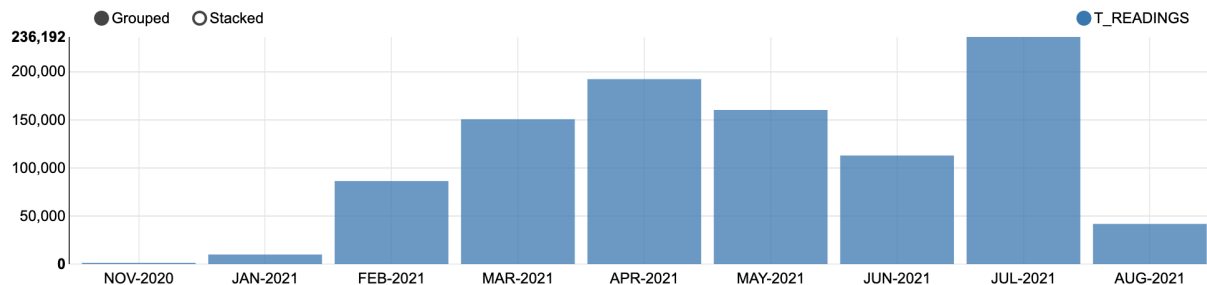


(a) Map of regions where case-study highways are located.

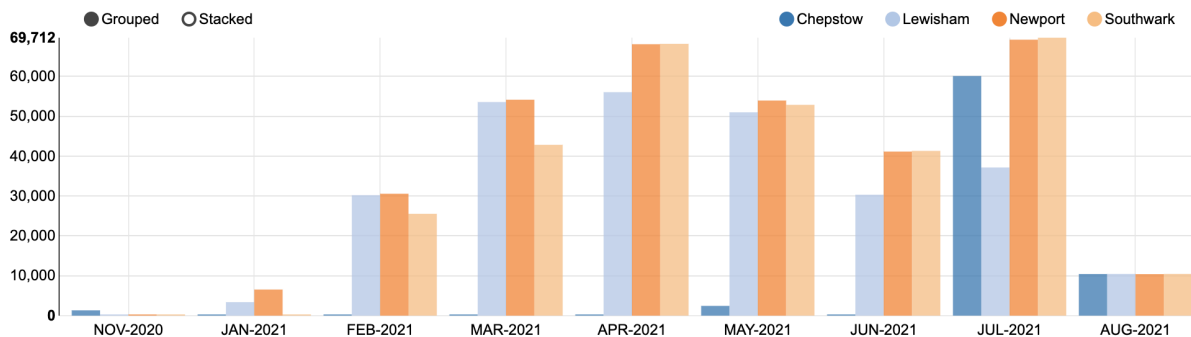


(b) Sensor distributions across the highways

Figure 4: Maps showing the distribution of 14 REVIS devices in four regions across the UK - Newport(1), Chepstow(1), Lewisham (6), Southwark(6). The devices in London were deployed to capture readings from the A302, A2209 and A1203 highways, while the devices in Newport and Chepstow were deployed close to the M4 and A48 highways, respectively



(a) Average monthly readings captured by deployed REVIS devices.



(b) This plot illustrates the number of readings captured per region.

Figure 5: Total monthly readings captured by deployed sensors between November 2020 and August 2021. These plots illustrate the amount of missing data in the first two months when some devices were offline. Chepstow had the lowest monitored readings overall.

226 4.2.1. The Impact of Weather on $PM_{2.5}$, PM_{10} and NO_2

227 Weather parameters influence the dispersion rates of pollutants (Barrera-Animas et al.
 228 2022). It is worthwhile to first check the correlation between the weather parameters before
 229 investigating the impact of weather on highway pollution. Figure 6 illustrates a correlation
 230 matrix constructed to identify the hierarchical similarities between these parameters which
 231 revealed a strong correlation between temperature, wind speed and wind direction. To un-
 232 derstand the effects of temperature on the four pollutants, the seasonal trends were plotted
 233 as shown in Figure 7. The average temperature for all four regions ranged between 8.6
 234 and 12.56°C in Winter, 9.73 and 19.76°C in spring and 19.41 and 21.78°C in summer. A
 235 regression analysis of temperature against each pollutant as presented in Table 2a depicts
 236 a positive correlation between $PM_{2.5}$ and PM_{10} and temperature in Newport, Southwark
 237 and Lewisham for spring and summer seasons. Chepstow had no correlation calculated for
 238 winter when there was no temperature reading recorded and a negative correlation in spring
 239 and summer. NO_2 had a negative correlation with temperature in all of these regions in
 240 winter and spring but had a positive correlation in Southwark and Lewisham in Summer.
 241 These findings corroborate studies that suggest that concentration levels are highest when
 242 the temperature is high (Pearce et al. 2011, Analitis et al. 2014).



Figure 6: Distance matrix of weather parameters using Pearson's correlation. A strong correlation can be noticed between “temp”, “temp_min”, “temp_max”, “wind_speed”, “wind_degree” and “feels_like”. There is also a discernible correlation between “clouds_all” and “humidity”/“windspeed”.

Table 2: Regression analysis of weather parameters vs pollutant concentration

Regions	Winter				Spring				Summer			
	temp($^{\circ}C$)	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$	temp($^{\circ}C$)	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$	temp($^{\circ}C$)	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$
Newport	8.60	0.53	0.03	0.46	9.73	0.56	0.59	0.48	19.58	0.32	0.61	0.51
Southwark	12.68	0.40	0.10	0.32	10.44	0.33	0.23	0.18	19.41	0.11	0.24	0.26
Lewisham	12.56	0.46	0.13	0	11.80	0.41	0.38	0.10	20.77	0.33	0.35	0.09
Chepstow	-	-	-	-	19.76	0.44	0.38	0.33	21.78	0.19	0.20	0.17

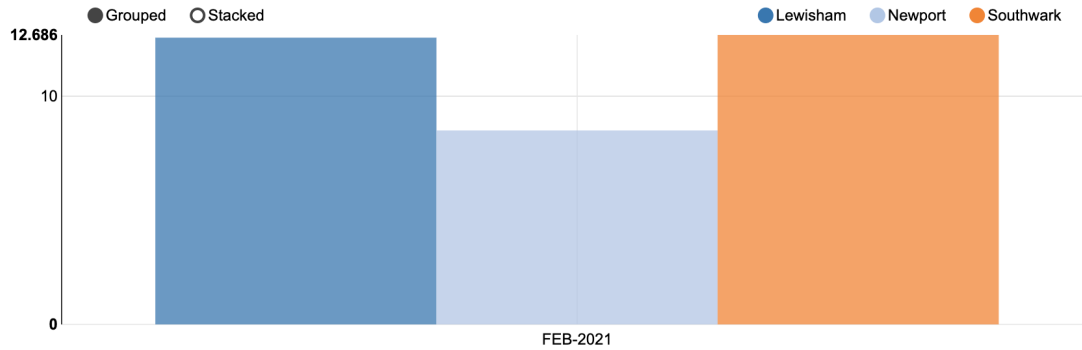
(a) Correlation between regional temperature and pollutants in spring, winter and summer

Regions	Winter				Spring				Summer			
	pressure	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$	pressure	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$	pressure	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$
Newport	1014	-0.10	0.42	0.38	1022.50	-0.06	0.12	0.31	1012.90	-0.13	0.33	0.55
Southwark	1018.50	-0.22	0.10	0.13	1026.10	-0.15	0.17	0.22	1015.30	-0.26	0.08	0.03
Lewisham	1018.80	-0.07	0.03	0.11	1026.30	-0.01	0.10	0.18	1014.20	-0.19	0.16	0.15
Chepstow	-	-	-	-	1009.50	0.44	-0.22	0.15	1007.20	0.19	0.10	0.09

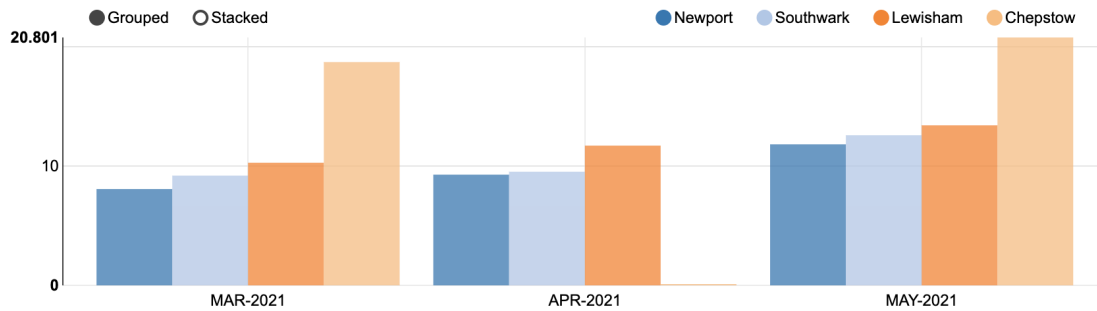
(b) Correlation between regional pressure and pollutants in spring, winter and summer

Regions	Winter				Spring				Summer			
	humidity($\%$)	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$	humidity($\%$)	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$	humidity($\%$)	$NO_2(r^2)$	$PM_{2.5}(r^2)$	$PM_{10}(r^2)$
Newport	90.96	0	-21	-18	73.54	2	-11	-3.40	70.85	1.30	-13.70	-4.80
Southwark	66.27	7	-1	-15	65.85	13	-6.50	-8.90	73.30	6.80	-3	-2.20
Lewisham	72.93	3	-8	-5	64.04	11	-15.20	-4.20	70.98	17.6	-17	-5.60
Chepstow	-	-	-	-	53.95	8	-1.70	-6	64.95	13.30	-3.4	-11.20

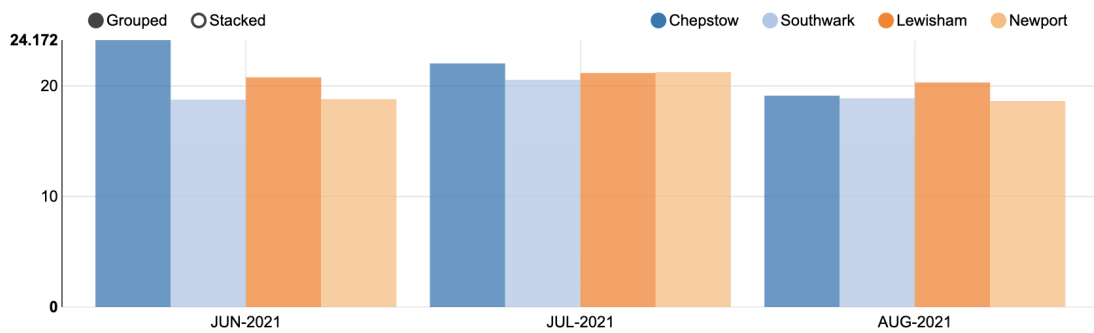
(c) Correlation between regional humidity and pollutants in spring, winter and summer



(a) Average winter temperature for all four regions

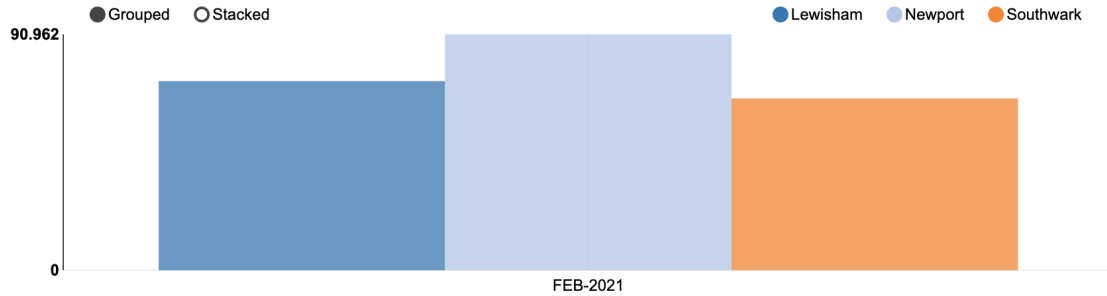


(b) Average spring temperature for all four regions

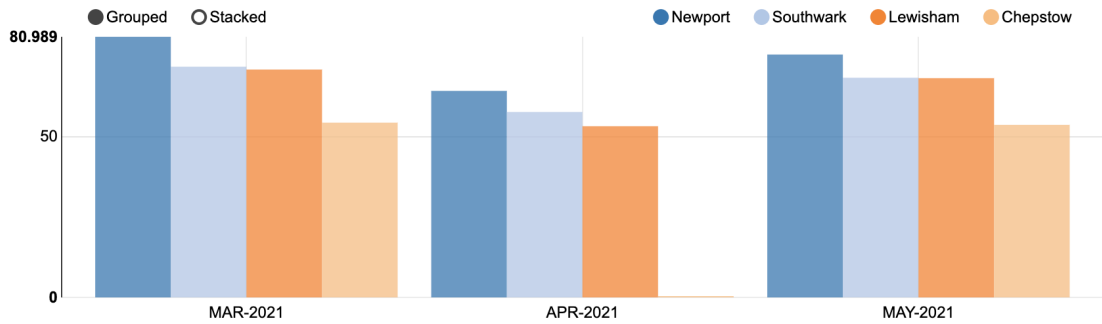


(c) Average summer temperature for all four regions

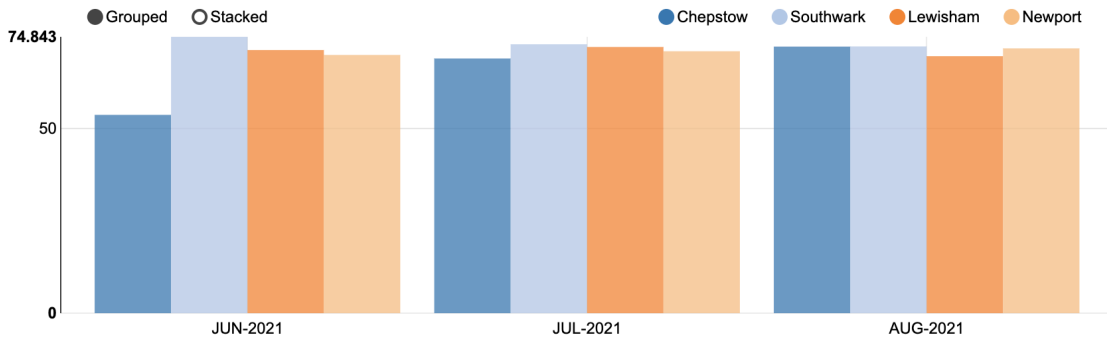
Figure 7: The seasonal trends for temperature in Newport, Southwark, Lewisham and Chepstow. Newport has the lowest temperature of 8.6°C in winter as there was also no reading recorded for Chepstow, as illustrated in plot (a). Chepstow had the highest average temperature of 19.76°C in spring and 21.78°C in summer, as shown in plots (b) and (c)



(a) Average winter humidity for all four regions



(b) Average spring humidity for all four regions



(c) Average summer humidity for all four regions

Figure 8: The seasonal trends for humidity in Newport, Southwark, Lewisham and Chepstow. Similar to temperature and pressure, no reading was captured for Chepstow in winter. However, the region recorded the least humidity of 53.95% in spring, as illustrated in plot (b). Newport had the highest average humidity of 90.96% in winter and 73.54% in spring, as shown in plots (a) and (b)

For pressure, the lowest readings were recorded in Chepstow during Spring and Summer seasons while Lewisham and Southwark recorded the highest pressures in spring. Table 2b summarises the pressure readings during these seasons and the correlation figures with the pollutants. The $PM_{2.5}$ and PM_{10} concentrations in Newport and Chepstow were positively correlated with pressure, indicating that an increase in atmospheric pressure will increase the concentration levels of these highway pollutants. All three pollutants negatively correlate with pressure in Southwark and Lewisham in spring but positive in winter and summer. The conclusion drawn from this result is a strong correlation between pressure and $PM_{2.5}$ and PM_{10} but a significant negative correlation with NO_2 . Figure 8 illustrates the average seasonal humidity across the regions with the lowest humidity value was recorded in Chepstow during summer and the highest in Newport during winter. It can be deduced from Table 2c that the three pollutants were negatively correlated with humidity for winter, spring and summer seasons. In particular, particulate matter ($PM_{2.5}$ and PM_{10}) are prone to be absorbed in the atmosphere as humidity increases. Naturally, rain results in higher relative humidity and soaks up these particles, resulting in a lower level of particulate in winter. (Odat 2009).

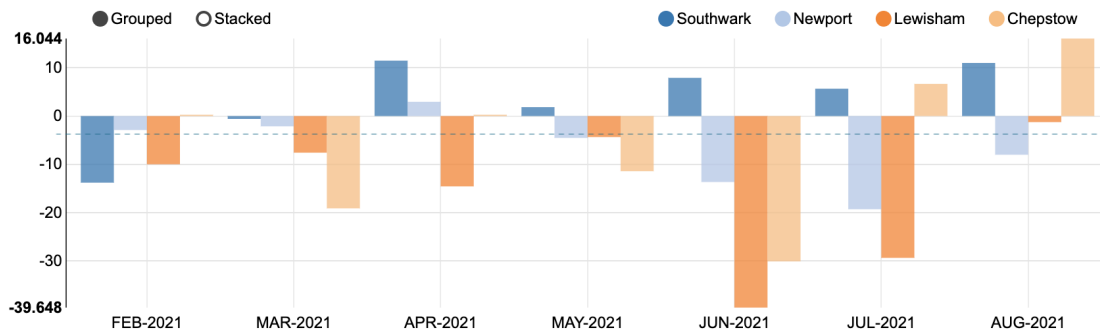
4.2.2. The Impact of Region on $PM_{2.5}$, PM_{10} and NO_2

Each region has its unique attributes which can influence the concentration level of pollutants measured over the experimentation period. Aside from the weather, other attributes such as the highway gradient, region terrain, residential development, background coefficient and traffic flow can also contribute to the concentration levels across regions. Although some of these attributes were not captured in this research, their effects on the captured concentration levels remain to be seen. This section presents some primary insights across the four regions in the dataset.

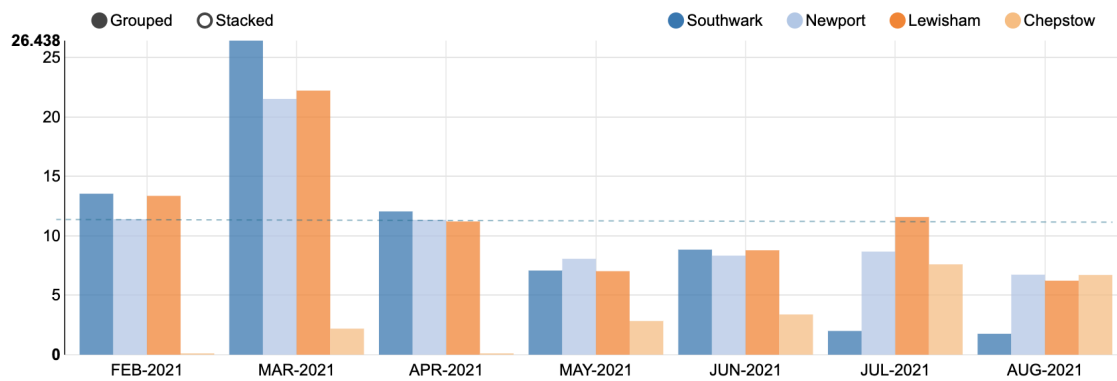
Table 3: Pollutant summary statistics based on region

Regions	NO_2				$PM_{2.5}$				PM_{10}			
	count	mean	min	max	count	mean	min	max	count	mean	min	max
Newport	40326	-6.85	-602.37	111.99	40326	11.41	0.28	745.45	40326	12.49	0.28	746.04
Southwark	38757	4.35	-714.97	1094.81	38757	10.27	0.55	4384.20	38757	11.35	0.60	6888.54
Lewisham	32986	-15.168	-1406.17	93.06	32986	12.42	0.60	277.02	32986	13.98	0.60	424.42
Chepstow	9138	7.33	-190.18	180.30	9138	7.31	0.43	127.45	9138	12.11	0.431	179.02

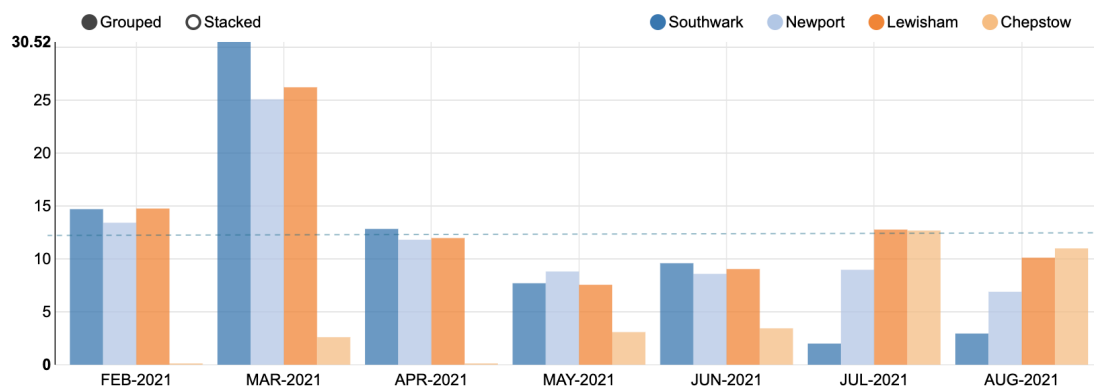
Table 3 shows that the average concentration levels across regions vary significantly, and this can be linked to the calibration accuracy of the sensor devices. Chepstow and Southwark seemed to have the most practical NO_2 averages, with Southwark having the highest. Lewisham has the highest $PM_{2.5}$ and PM_{10} average of $12.42\mu g/m^3$ and $13.98\mu g/m^3$, respectively. This analysis and the plot in Figure 9 reveal some prevalent calibration issues within the recorded values, which were sometimes exaggerated, as in the case of the maximum values for $PM_{2.5}$ and PM_{10} . Nevertheless, a one-way ANOVA variance test carried out to check the variance in NO_2 , $PM_{2.5}$ and PM_{10} by region resulted in p values of $2.36e-4$, $1.45e-3$



(a) Monthly NO_2 average for all four regions



(b) Monthly $PM_{2.5}$ average for all four regions



(c) Monthly PM_{10} average for all four regions

Figure 9: Plots highlighting the varying monthly averages for the three monitored pollutants. These averages varied significantly and are an indication that some influential factors may have affected the concentration levels

275 and 1.68e-4, respectively. This result indicates that the impact of regions on the concentration
276 levels of these three pollutants is notable.

277 4.3. Forecasting Model Training and Evaluation

278 Fastai was used for data pre-processing and model training. The library is built on the
279 PyTorch framework and allows quick analysis using its readily encoded best practices. The
280 aim was to develop a model capable of efficiently making hourly predictions of the pollutant
281 of interest. This section introduces the data processing procedure, the network’s architecture
282 used for training and the validation method.

283 4.3.1. Meteorology Data Integration and Dataset Pre-Processing

284 Weather data such as wind speed and direction, precipitation, visibility, pressure, cloud
285 cover, dew point, and wind gust which were not captured by the REVIS devices, were inte-
286 grated from OpenWeather. Also, ozone data from the AURN stations were integrated into
287 the dataset to be analysed and used for training. These integration exemplify the integration
288 capabilities of the framework while enriching the data needed to train an estimation model.
289 Appendix A presents a complete list of the columns, their description and data types be-
290 fore processing. An SQL procedure for automatically generating SQL codes such as the one
291 illustrated in Figure 10 was implemented to summarise the pollution data. This generated
292 hourly, 3-hourly and 6-hourly summaries of the pollutant concentration levels with the aim
293 of capturing periodicity within the training data.

294 Three key data pre-processors: *categorify*, *fillMissing* and *normalize* from fastai were
295 adopted for additional data pre-processing. These pre-processors map categorical columns to
296 distinct categories, replaces null values with column median values and normalises continuous
297 columns by subtracting the mean and dividing by the standard deviation. The “*add_datepart*”
298 helper function of the library allows the specification of the date column which generates
299 additional predictors such as “*Year*”, “*DayofWeek*”, “*DayOfYear*”, “*Is_Month_End*” and so
300 on. Appendix B highlights the list of categorical and continuous variables in the dataset after
301 processing.

302 4.3.2. Validation Set Creation and Training Architecture

303 Model training is typically initiated by splitting the dataset into training, test and val-
304 idation datasets. As the name implies, training data is used for training, while validation
305 data is used for selecting the model that works best after verification using the test data.
306 It is customary to randomise the dataset before splitting when there is a class imbalance
307 - stratification; but since this problem is similar to a time-series problem where the date
308 order is important, the validation and test sets can not be randomly selected. The common
309 practice is to select the last few weeks or months of the dataset for validation and testing.
310 Our dataset of 991662 rows and 34 columns had no class imbalance for the three pollutants
311 which meant that stratification was not necessary. An experiment with different numbers of
312 last days was carried out to determine the best validation approach, and the last 45 days of
313 the dataset from July and August were eventually chosen for validation(15 days) and test(30
314 days). Fastai’s “TrainTestSplitter” class was used to implement this division.

```

SELECT city_name, lat, lon, TO_CHAR(edate,'yyyy-mm-dd hh24:mi:ss') edate, Rain_desc, Rain_1h, Rain_3h, Snow_1h, Snow_3h,
Drizzle_desc, Fog_desc, Clouds_desc, Haze_desc, Mist_desc, Clear_desc, Snow_desc, Thunderstorm_desc, temp, temp_min,
temp_max, feels_like, pressure, humidity, wind_speed, wind_deg, clouds_all,
ROUND(Ozone, 4) Ozone, ROUND(AVG(Ozone) OVER (ORDER BY edate ROWS BETWEEN 6 PRECEDING AND 1 PRECEDING), 4) Ozone_avg6h, Ozone_factor,
ROUND(no, 4) no, ROUND(AVG(no) OVER (ORDER BY edate ROWS BETWEEN 6 PRECEDING AND 1 PRECEDING), 4) no_avg6h, no_factor,
ROUND(no2, 4) no2, ROUND(AVG(no2) OVER (ORDER BY edate ROWS BETWEEN 6 PRECEDING AND 1 PRECEDING), 4) no2_avg6h, no2_factor,
ROUND(nono2, 4) nono2, ROUND(AVG(nono2) OVER (ORDER BY edate ROWS BETWEEN 6 PRECEDING AND 1 PRECEDING), 4) nono2_avg6h, nono2_factor,
ROUND(pm10, 4) pm10, ROUND(AVG(pm10) OVER (ORDER BY edate ROWS BETWEEN 6 PRECEDING AND 1 PRECEDING), 4) pm10_avg6h, pm10_factor,
ROUND(pm25, 4) pm25, ROUND(AVG(pm25) OVER (ORDER BY edate ROWS BETWEEN 6 PRECEDING AND 1 PRECEDING), 4) pm25_avg6h, pm25_factor
FROM (SELECT city_name, lat, lon, edate, TO_CHAR(edate,'mm-dd hh24') edate_hr, Rain_desc, Rain_1h, Rain_3h,
Snow_1h, Snow_3h, Drizzle_desc, Fog_desc, Clouds_desc, Haze_desc,
Mist_desc, Clear_desc, Snow_desc, Thunderstorm_desc, temp, temp_min, temp_max, feels_like,
pressure, humidity, wind_speed, wind_deg, clouds_all,
nvl(last_value(nullif((CASE WHEN Ozone<0 THEN null ELSE Ozone END), 0)) IGNORE NULLS OVER (ORDER BY edate), 0) Ozone,
nvl(last_value(nullif((CASE WHEN no<0 THEN null ELSE no END), 0)) IGNORE NULLS OVER (ORDER BY edate), 0) no,
nvl(last_value(nullif((CASE WHEN no2<0 THEN null ELSE no2 END), 0)) IGNORE NULLS OVER (ORDER BY edate), 0) no2,
nvl(last_value(nullif((CASE WHEN nono2<0 THEN null ELSE nono2 END), 0)) IGNORE NULLS OVER (ORDER BY edate), 0) nono2,
nvl(last_value(nullif((CASE WHEN pm10<0 THEN null ELSE pm10 END), 0)) IGNORE NULLS OVER (ORDER BY edate), 0) pm10,
nvl(last_value(nullif((CASE WHEN pm25<0 THEN null ELSE pm25 END), 0)) IGNORE NULLS OVER (ORDER BY edate), 0) pm25
FROM emissions_main) JOIN emission_factors USING (city_name, edate_hr);

```

Figure 10: Auto-SQL generation to pre-process the dataset. An SQL command which generates 3-hour and 6-hour pollutant averages from the preceding readings is depicted.

315 Suitable optimisers, loss functions and activation functions had to be selected from an
316 array of available options. Series of experimentation were carried out on popular optimisation
317 functions such as *SGD*, *RMSProp*, *LAMB*, *LARS* and *Adam* and regression loss functions like
318 *BCELossFlat*, *MSELossFlat* and *L1LossFlat* before deciding the most suitable. Eventually,
319 *Adam* optimiser and *MSELossFlat* were chosen for model training. *Bayesian-optimization*
320 library was used to test and optimise the number of architecture layers, the size of each layer
321 and dropout rates for the network. The final architecture used to train the model was made
322 up of 14 embedding layers, 3 dropout layers, 3 batchnorm1d layers, 3 linear layers and 2
323 ReLU activation functions. The embedding layer was adopted for improved performance as
324 inspired by the architecture proposed in Guo & Berkahn (2016). Finally, the learning rate
325 finder (*lr_finder*) function of *TabularLearner* class was used to determine the best learning rate
326 to be used for training. This resulted in a minimum value of $2.5e^{-4}$, and steep value of $1.3e^{-4}$.
327 Figure 11 below shows the plot of the learning rate against the loss. Experts recommend
328 selecting the learning rate at the point where the plot starts to dip. (i.e., 10^{-4}).

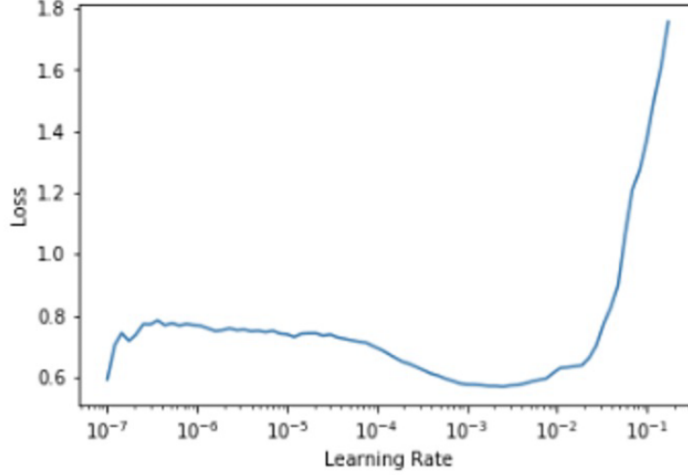


Figure 11: The model’s training loss against the learning rate to determine the appropriate learning rate. The learning rate was fixed at the point where the plot started dipping (i.e., 10^{-4})

329 4.3.3. Model Evaluation

330 In this section, the results of the deep learning model developed are presented. The model
 331 was trained to make day-ahead predictions of the three pollutants, but first, an appropriate
 332 evaluation metric had to be selected. The top metrics for regression problems are mean
 333 squared error/root mean squared error(MSE/RMSE), mean absolute error(MAE) and R
 334 Square. The fastai library has two variants of RMSE: *rmse* and *exp_rmse*. The mean absolute
 335 error and root mean squared error (*exp_rmse* variant), defined as shown in equations 2 and
 336 3 below, were selected as the metrics for evaluating the developed model.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (2)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|^2} \quad (3)$$

337 Figure 12 illustrates the model training and validation losses after 20000 epochs. It is
 338 noteworthy that the training loss gradually as the number of epoch increased. The validation
 339 loss took a slightly different pattern and dropped significantly after 2500 epochs but became
 340 steady for the remaining training epochs. The final MAE and exponential RMSE after
 341 training were 0.350 and 1.591 respectively. Figure 13 captures the actual NO_2 concentration
 342 levels (highlighted in blue) and the model’s day ahead prediction(highlighted in red). The
 343 difference in the model’s predicted NO_2 and actual values is slight, and the predicted values
 344 were close to the actual.

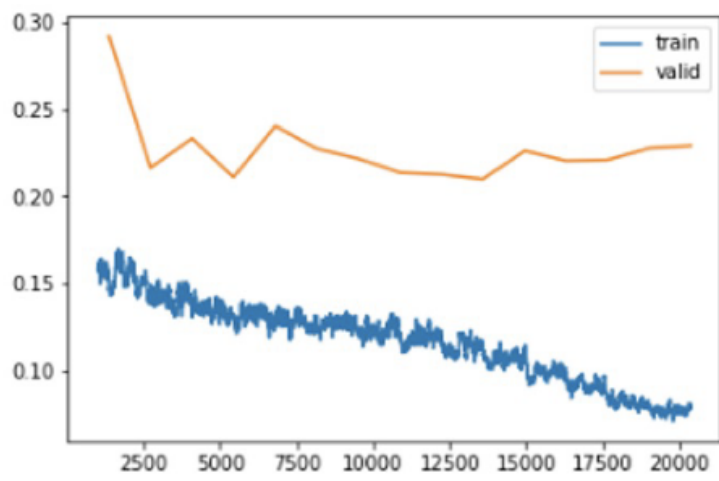


Figure 12: A plot showing the model’s training and validation losses against the number of epochs. It is worth noting that there was a gradual decrease in both losses as the training epochs increased which indicates that the model was learning. Further training beyond 20000 epochs would have either resulted in overfitting or no further drop in both losses

```
learn.show_results()
```

R	NO2_AVG6H	NO2_FACTOR	NONO2_AVG6H	NONO2_FACTOR	PM10_AVG6H	PM10_FACTOR	PM25_AVG6H	PM25_FACTOR	Elapsed	NO2	NO2_pred
4	-0.885981	-0.332012	-0.567841	0.662468	-0.162312	-0.699817	-0.219881	-0.814794	1.733253	2.034706	2.129269
6	-0.748497	0.903192	-0.447542	0.058060	-0.162312	0.790438	-0.219881	1.451507	1.800624	2.850389	3.122491
8	-1.236534	-0.112674	-0.615339	-0.334880	-0.162312	0.790438	-0.219881	0.883862	1.835902	1.898369	2.188342
7	-0.354666	2.035958	-0.354947	0.910231	-0.162312	1.176936	-0.219881	1.997850	1.832319	2.943470	3.006840
1	0.368544	0.594365	-0.066416	0.100828	-0.162312	1.176936	-0.219881	1.284325	1.764231	2.849961	3.192918
4	-1.124172	-0.726089	-0.597230	-0.589171	-0.162312	0.790438	-0.219881	0.728293	1.843501	2.016155	2.143995
9	0.467780	0.327762	1.666261	0.880060	-0.162312	-0.274668	-0.219881	-0.643668	1.757341	3.549833	3.641593
5	-0.125720	0.678973	-0.298762	0.533345	-0.162312	0.790438	-0.219881	0.549499	1.785651	2.364968	2.505113
9	0.341336	2.563796	-0.033054	0.957226	-0.162312	1.541260	-0.219881	1.607077	1.766261	3.858867	3.541561

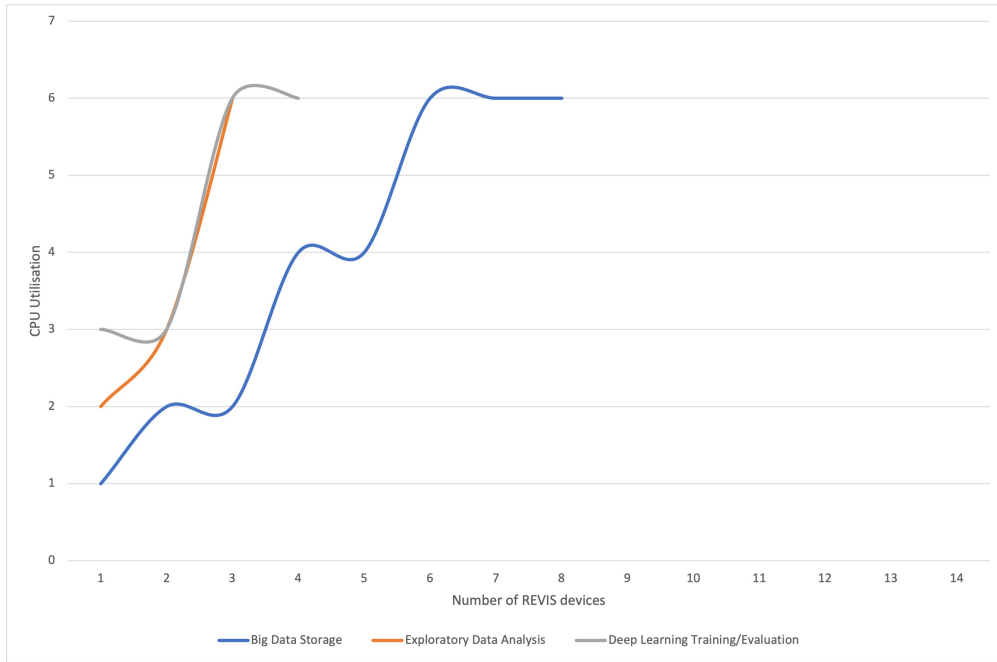
Figure 13: An illustration of captured NO_2 pollutant readings (blue highlight) and the deep learning model predictions (red highlight). These results were derived from an evaluation using the validation dataset. It should be pointed out that the model’s predictions are not too far off the actual readings.

345 **4.4. Evaluating the Scalability Performance of the REVIS System**

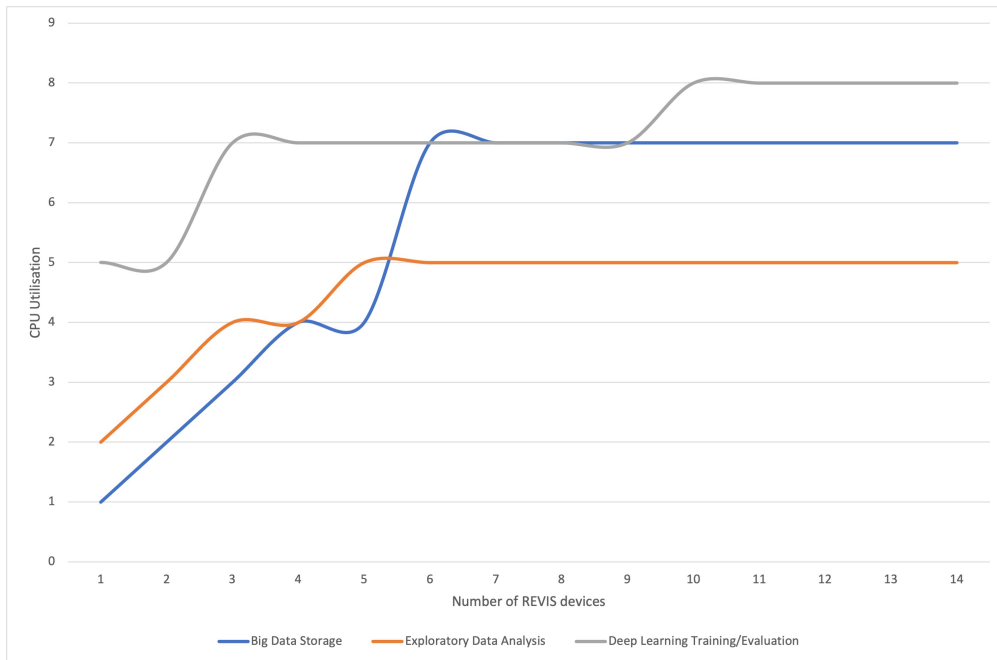
346 The REVIS system was tested for scalability using the IoT asset monitoring tool and
 347 database performance hub of two different oracle cloud instances. The fourteen REVIS
 348 devices were deployed sequentially to capture both system’s response time and throughput.
 349 The first experiment was run on a bare metal cloud instance with specifications as shown in
 350 Table 4. Figure 14a shows the performance of this cloud instance as it could not scale past
 351 8 devices and exploded at 3 and 4 devices for EDA and deep learning analysis. However,
 352 the GPU cloud instance performed better due to its auto-scale feature. Figure 14b shows a
 353 plot of the CPU cores utilised for exploratory data analysis, data storage and deep learning
 354 analysis as the number of deployed devices increased. It can be observed that the number of
 355 CPU cores increased gradually for each task and then stabilised at some point. The system
 356 was able to scale up its resources according to the computation/storage requirements. For
 357 the database performance, the test was run between November 2020 and Jan 2021 on the
 358 GPU instance and evaluated for utilisation, execution count, number of running statements
 359 and number of sessions metrics as shown in Figure 15. The maximum GPU utilisation was
 360 under 20% even with over 1.5 million execution queries.

Table 4: Hardware specifications of the two oracle cloud instances used to test scalability

Name	Instance Type	Processor	GPU type	CPU cores	CPU memory	GPU memory
Compute - Ampere A1 - OCPU	Bare Metal	OCPU	-	6	32GB	-
VM.GPU2.1	GPU	Pascal	1 NVIDIA P100	12	72GB	16GB

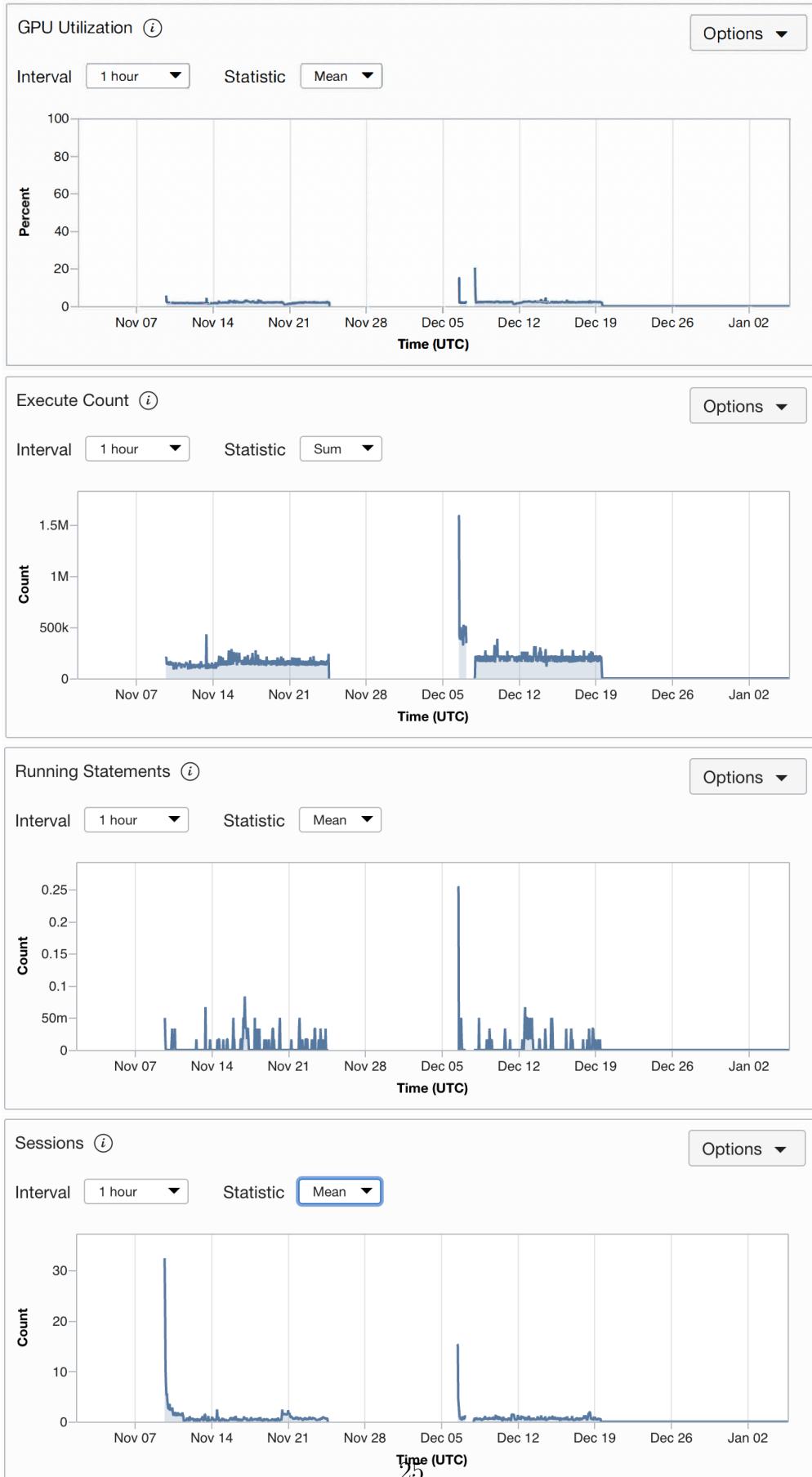


(a) System performance of the bare metal instance as the number of REVIS devices increased.



(b) System performance of the GPU instance as the number of REVIS devices increased.

Figure 14: Plots of bare metal vs GPU instance as number of devices increased



(a) Metrics showing system's performance after the deployment of 14 sensors

Figure 15: Plots of scalability metrics showing database performance as the number of devices increased

5. Discussion and Implication of Study

The REVIS system was used to demonstrate the possibility of optimising the cost, efficiency and environmental impact of hardware IoT devices through the development, calibration and deployment of monitoring units to capture real-time pollution data on highways. The devices were developed through an excellent design of both analogue and digital circuitry around it and an iterative approach of calibration and performance optimisation. Although we were able to address the energy interference and cross-sensitivity issues of existing sensing devices, the developed units still had some NO_2 data inconsistencies which were directly linked to the chosen pollutant sensor. A probable solution is the adoption of machine learning techniques for sensor data calibration. This technique is increasingly becoming popular and has been explored in studies such as that of Zimmerman et al. (2018) and Si et al. (2020). Nevertheless, our implementation still demonstrates the hardware layer of the proposed framework and also the effectiveness of carefully designed low-cost and environmentally-friendly sensors in capturing and processing accurate data on highway air quality.

It is important to note that sensors data alone are not sufficient for ensuring accuracy in air quality forecasting models. There are a number of air quality data sources, which exist separately but can provide better insights about air quality if well explored and integrated. An important aspect of this study is to integrate missing or inaccurate data from heterogeneous sources to enhance forecasting accuracy of the developed deep learning model. The essence of this layer is to ensure that data not captured in the hardware layer by the monitoring devices can be integrated into the system to improve the performance. Similarly, an exploratory analysis on the captured and integrated data was conducted to evaluate the impact of different parameters on pollutant concentration. It is well established in literature that weather parameters such as rainfall and temperature influence the dispersion rates of pollutants (Barrera-Animas et al. 2022). Hence, there is need for a more coordinated approach such as the one proposed in this study to manage multiple data sources, which are relevant for accurately forecasting air quality on highways through common data environment and data integration.

Finally, this study set out to develop and evaluate a baseline deep learning model to make hourly predictions of $PM_{2.5}$, PM_{10} and NO_2 concentration levels. The problem is modelled as a structured data with additional features added to extend a typical time-series problem. This method allows us to explore entity embeddings for the categorical features. The performance of the baseline model on individual pollutant concentration is good, thereby suggesting that this approach of modelling is practicable. An improved forecasting model is directly applicable to highways where air quality sensing devices are not available but data on other features such as traffic flow and weather are captured. These models can then be used as a substitute to estimate the air quality on these highways. The scalability test performed on the REVIS system indicated it was able to scale up its resources according to the computation/storage requirements. This study has addressed an important aspect of air quality management on highways, provided a scalable solution for academics and industry practitioners; and a pathway for policy makers and highway regulators to make more informed decisions.

403 6. Conclusion

404 A cost-effective deep learning framework for ubiquitous monitoring and predicting pollu-
405 tant concentration levels on UK highways was proposed in this study. An implementation of
406 the framework was demonstrated using the REVIS system. Details of the development of the
407 REVIS IoT hardware for data collection, the configuration of big data tools for data storage
408 and the and results of trained deep learning forecasting models were reported. The scalability
409 feature of the framework was also highlighted using two cloud instances with different com-
410 putational resources. This study showed that real-time monitoring and forecasting could be
411 achieved with the right computational resources. Although the scope of the research was lim-
412 ited to NO_2 $PM_{2.5}$ and PM_{10} pollutants and evaluation using deep learning, future research
413 can focus on investigating other pollutants such as CO_2 , SO_2 and Ozone as well as other
414 machine learning approaches for estimation. This study is a part of a series of publications
415 highlighting research findings on the REVIS project. This is just a scratch on the surface
416 of our research outputs, as other articles will elaborate on developing more advanced deep
417 learning models. Future research will integrate other relevant highway attributes such as
418 traffic flow, highway terrain, background concentration, and pollutant characteristics such as
419 washout coefficient, dispersion rate, and emission, critical to developing highway air quality
420 estimation models.

421 Acknowledgement

422 The authors would like to express their sincere gratitude to InnovateUK (Grant Applica-
423 tion No 10137 and File No 104367) and EPSRC(Grant Ref No EP/N509012/1) for providing
424 the financial support for this study.

425 **Appendix A: Data summary for pollutant estimation before processing**

S/No	Column	Column Description	Non-Null Count	Data type
1	city_name	The name of the city of interest	991662 non-null	object
2	lat	The geographic coordinate of the city of interest (Latitude)	991662 non-null	float64
3	lon	The geographic coordinate of the city of interest (Longitude)	991662 non-null	float64
4	date	The observation time to include date, time, hour and second	991662 non-null	datetime64[ns]
5	rain_desc	Description of measured precipitation	5975 non-null	object
6	rain_1h	Integrated average hourly precipitation measurement (mm)	5658 non-null	float64
7	rain_3h	Integrated precipitation measurement averaged over 3 hrs preceding the observation time (mm)	65 non-null	float64
8	snow_1h	Integrated average hourly snow depth measurement (cm)	77 non-null	float64
9	snow_3h	Integrated snow depth measurement averaged over 3 hrs preceding the observation time (cm)	4 non-null	float64
10	drizzle_desc	Description of measured drizzle	244 non-null	object
11	fog_desc	Description of measured fog	193 non-null	object
12	clouds_desc	Description of measured clouds	72395 non-null	object
13	haze_desc	Description of measured haze	46 non-null	object
14	mist_desc	Description of measured mist	312 non-null	object
15	clear_desc	Description of measured clear	11342 non-null	object
16	snow_desc	Description of measured snow	103 non-null	object
17	storm_desc	Description of measured thunderstorm	1 non-null	object
18	temp	Captured average hourly temperature (°C)	991662 non-null	float64
19	temp_min	Captured minimum temperature over a 24-hr period (°C)	991662 non-null	float64
20	temp_max	Captured maximum temperature over a 24-hr period (°C)	991662 non-null	float64

21	feels_like	Integrated measurement of human impression of weather (K)	991662 non-null	float64
22	pressure	Captured average hourly pressure (hPa)	991662 non-null	int64
23	humidity	Captured average hourly relative humidity (ϕ)	991662 non-null	int64
24	wind_speed	Integrated average hourly wind speed (knots)	991662 non-null	float64
25	wind_direction	Integrated average hourly wind direction (true degrees)	991662 non-null	int64
26	clouds_all	Integrated hourly measurement of cloudiness (%)	991662 non-null	float64
27	ozone	Integrated average hourly ozone ($\mu g/m^3$)	181233 non-null	float64
28	ozone_avg6h	Integrated ozone readings averaged over 6 hrs preceding the observation time ($\mu g/m^3$)	181233 non-null	float64
29	NO_2	Captured average hourly NO_2 (ppb)	121207 non-null	float64
30	NO_2 _avg6h	Captured NO_2 readings averaged over 6 hrs preceding the observation time (ppb)	121207 non-null	float64
31	PM_{10}	Captured average hourly PM_{10} ($\mu g/m^3$)	121207 non-null	float64
32	PM_{10} _avg6h	Captured PM_{10} readings averaged over 6 hrs preceding the observation time ($\mu g/m^3$)	121207 non-null	float64
33	$PM_{2.5}$	Captured average hourly $PM_{2.5}$ ($\mu g/m^3$)	121207 non-null	float64
34	$PM_{2.5}$ _avg6h	Captured $PM_{2.5}$ readings averaged over 6 hrs preceding the observation time ($\mu g/m^3$)	121207 non-null	float64

426

427 **Appendix B: List of attributes after processing and classification as categorical**
428 **or continuous**

S/No	Attribute Name	Attribute Type
1	city_name	Categorical
2	lat	Categorical
3	lon	Categorical
4	year	Categorical
5	month	Categorical

6	week	Categorical
7	day	Categorical
8	dayofweek	Categorical
9	dayofyear	Categorical
10	is_month_end	Categorical
11	is_month_start	Categorical
12	is_quarter_end	Categorical
13	is_quarter_start	Categorical
14	is_year_end	Categorical
15	is_year_start	Categorical
16	rain_1h	Continuous
17	snow_1h	Continuous
18	temp	Continuous
19	temp_min	Continuous
20	temp_max	Continuous
21	feels_like	Continuous
22	pressure	Continuous
23	humidity	Continuous
24	wind_speed	Continuous
25	wind_direction	Continuous
26	clouds_all	Continuous
27	ozone	Continuous
28	ozone_avg6h	Continuous
29	<i>no₂</i>	Continuous
30	<i>no₂_avg6h</i>	Continuous
31	<i>pm_{2.5}</i>	Continuous
32	<i>pm_{2.5}_avg6h</i>	Continuous

34	PM_{10}	Continuous
34	PM_{10_avg6h}	Continuous

430 References

- 431 Ahmed, E., Ahmed, A., Yaqoob, I., Shuja, J., Gani, A., Imran, M. & Shoaib, M. (2017),
432 ‘Bringing computation closer towards user network: Is edge computing the solution?’,
433 *IEEE Communications Magazine* **55**, 138 – 144.
- 434 Akinosho, T. D., Oyedele, L. O., Bilal, M., Ajayi, A. O., Delgado, M. D., Akinade, O. O.
435 & Ahmed, A. A. (2020), ‘Deep learning in the construction industry: A review of present
436 status and future innovations’, *Journal of Building Engineering* p. 101827.
- 437 Alléon, A., Jauvion, G., Quennehen, B. & Lissmyr, D. (2020), ‘Plumenet: Large-scale air
438 quality forecasting using a convolutional lstm network’, *arXiv preprint arXiv:2006.09204* .
- 439 Alvanchi, A., Rahimi, M., Mousavi, M. & Alikhani, H. (2020), ‘Construction schedule, an
440 influential factor on air pollution in urban infrastructure projects’, *Journal of Cleaner
441 Production* **255**, 120222.
- 442 Analitis, A., Michelozzi, P., D’Ippoliti, D., De’Donato, F., Menne, B., Matthies, F., Atkinson,
443 R. W., Iñiguez, C., Basagaña, X., Schneider, A. et al. (2014), ‘Effects of heat waves on
444 mortality: effect modification and confounding by air pollutants’, *Epidemiology* pp. 15–22.
- 445 Badura, M., Batog, P., Drzeniecka-Osiadacz, A. & Modzel, P. (2018), Optical particulate
446 matter sensors in pm2. 5 measurements in atmospheric air, in ‘E3S Web of Conferences’,
447 Vol. 44, EDP Sciences, p. 00006.
- 448 Barikayeva, N., Nikolenko, D. & Ivanova, J. (2018), About forecasting air pollution in the
449 construction of highways, in ‘IOP Conference Series: Materials Science and Engineering’,
450 Vol. 463, IOP Publishing, p. 042016.
- 451 Barrera-Animas, A. Y., Oyedele, L. O., Bilal, M., Akinosho, T. D., Delgado, J. M. D.
452 & Akanbi, L. A. (2022), ‘Rainfall prediction: A comparative analysis of modern ma-
453 chine learning algorithms for time-series forecasting’, *Machine Learning with Applications*
454 **7**, 100204.
- 455 Barthwal, A. & Acharya, D. (2018), An internet of things system for sensing, analysis &
456 forecasting urban air quality, in ‘2018 IEEE International Conference on Electronics, Com-
457 puting and Communication Technologies (CONECCT)’, IEEE, pp. 1–6.
- 458 Bilal, M. & Oyedele, L. O. (2020), ‘Guidelines for applied machine learning in construc-
459 tion industry—a case of profit margins estimation’, *Advanced Engineering Informatics*
460 **43**, 101013.
- 461 Borghi, F., Spinazzè, A., Campagnolo, D., Rovelli, S., Cattaneo, A. & Cavallo, D. M. (2018),
462 ‘Precision and accuracy of a direct-reading miniaturized monitor in pm2. 5 exposure as-
463 sessment’, *Sensors* **18**(9), 3089.
- 464 Budde, M., Müller, T., Laquai, B., Streibl, N., Schwarz, A., Schindler, G., Riedel, T., Beigl,
465 M. & Dittler, A. (2018), Suitability of the low-cost sds011 particle sensor for urban pm-
466 monitoring, in ‘3rd International Conference on Atmospheric Dust’.

- 467 Carullo, A., Corbellini, S. & Grassini, S. (2007), ‘A remotely controlled calibrator for chem-
468 ical pollutant measuring-units’, *IEEE Transactions on Instrumentation and Measurement*
469 **56**(4), 1212–1218.
- 470 Chen, J., Li, K., Deng, Q., Li, K. & Philip, S. Y. (2019), ‘Distributed deep learning model
471 for intelligent video surveillance systems with edge computing’, *IEEE Transactions on*
472 *Industrial Informatics* .
- 473 Chen, M., Wang, S. & Xu, Q. (2015), ‘Multiobjective optimization for air-quality monitoring
474 network design’, *Industrial & Engineering Chemistry Research* **54**(31), 7743–7750.
- DEFRA (2019), ‘Clean air strategy’.
URL: https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/air-strategy-2019.pdf
- 475 DEFRA (2020), ‘Air quality appraisal: impact pathways approach’.
476 **URL:** [https://www.gov.uk/government/publications/assess-the-impact-of-air-quality/air-](https://www.gov.uk/government/publications/assess-the-impact-of-air-quality/air-quality-appraisal-impactpathways-approach)
477 [quality-appraisal-impactpathways-approach](https://www.gov.uk/government/publications/assess-the-impact-of-air-quality/air-quality-appraisal-impactpathways-approach)
- 478 Guo, C. & Berkhahn, F. (2016), ‘Entity embeddings of categorical variables’, *arXiv preprint*
479 *arXiv:1604.06737* .
- 480 Karagulian, F., Barbieri, M., Kotsev, A., Spinelle, L., Gerboles, M., Lagler, F., Redon, N.,
481 Crunaire, S. & Borowiak, A. (2019), ‘Review of the performance of low-cost sensors for air
482 quality monitoring’, *Atmosphere* **10**(9), 506.
- 483 Karner, A. A., Eisinger, D. S. & Niemeier, D. A. (2010), ‘Near-roadway air quality: synthesiz-
484 ing the findings from real-world data’, *Environmental science & technology* **44**(14), 5334–
485 5344.
- 486 Lee, M., Lin, L., Chen, C.-Y., Tsao, Y., Yao, T.-H., Fei, M.-H. & Fang, S.-H. (2020),
487 ‘Forecasting air quality in taiwan by using machine learning’, *Scientific reports* **10**(1), 1–
488 13.
- 489 Mabahwi, N. A. B., Leh, O. L. H. & Omar, D. (2014), ‘Human health and wellbeing: Human
490 health effect of air pollution’, *Procedia-Social and Behavioral Sciences* **153**, 221–229.
- 491 Odat, S. (2009), ‘Diurnal and seasonal variation of air pollution at al-hashimeya town, jor-
492 dan’, *Earth Environ Sci* **2**, 1–6.
- 493 ONS (2018), ‘2017 uk greenhouse gas emissions, provisional figures’, *Statistical Release: Na-*
494 *tional Statistics* .
- 495 ONS (2020), ‘Population estimates for regions in england and wales by sex and age’, Avail-
496 able from: <https://www.statista.com/statistics/1064772/population-of-london-by-gender/>
497 [Accessed: 15-12-2021].
- 498 ONS (2021), ‘Labour market in the regions of the uk: October 2021’, Available
499 from: [https://www.gov.uk/government/statistics/labour-market-in-the-regions-of-the-uk-](https://www.gov.uk/government/statistics/labour-market-in-the-regions-of-the-uk-october-2021)
500 [october-2021](https://www.gov.uk/government/statistics/labour-market-in-the-regions-of-the-uk-october-2021) [Accessed: 15-12-2021].

501
502
503

504
505

506
507
508

509
510
511

512
513
514

515

516
517
518

519
520
521
522

523
524
525

526
527

528
529
530

531
532
533

Pearce, J. L., Beringer, J., Nicholls, N., Hyndman, R. J. & Tapper, N. J. (2011), ‘Quantifying the influence of local meteorology on air quality using generalized additive models’, *Atmospheric Environment* **45**(6), 1328–1336.

Public Health England (2019), ‘Review of interventions to improve outdoor air quality and public health’.

URL: https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/2019-2018572.pdf

Sergeev, A. & Del Balso, M. (2018), ‘Horovod: fast and easy distributed deep learning in tensorflow’, *arXiv preprint arXiv:1802.05799*.

Shrivastava, A., Gupta, V. B. et al. (2011), ‘Methods for the determination of limit of detection and limit of quantitation of the analytical methods’, *Chronicles of young scientists* **2**(1), 21.

Si, M., Xiong, Y., Du, S. & Du, K. (2020), ‘Evaluation and calibration of a low-cost particle sensor in ambient conditions using machine-learning methods’, *Atmospheric Measurement Techniques* **13**(4), 1693–1707.

Statswales (2020), ‘Summary statistics for wales, by region: 2020’, Available from: <https://gov.wales/sites/default/files/statistics-and-research/2020-05/summary-statistics-regions-wales-2020-629.pdf> [Accessed: 15-12-2021].

TFL (2019), ‘Travel in london: Understanding our diverse communities 2019’.

Umadevi, K. & Geraldine Bessie Amali, D. (2020), Data visualization and analysis for air quality monitoring using ibm watson iot platform, *in* ‘Data Visualization’, Springer, pp. 15–32.

Vohra, K., Marais, E. A., Suckra, S., Kramer, L., Bloss, W. J., Sahu, R., Gaur, A., Tripathi, S. N., Van Damme, M., Clarisse, L. et al. (2021), ‘Long-term trends in air quality in major cities in the uk and india: A view from space’, *Atmospheric Chemistry and Physics* **21**(8), 6275–6296.

World Bank (2022), ‘The global health cost of pm2.5 air pollution: A case for action beyond 2021’.

URL: <https://openknowledge.worldbank.org/handle/10986/36501>

Zhang, K. & Batterman, S. (2013), ‘Air pollution and health risks due to vehicle traffic’, *Science of the total Environment* **450**, 307–316.

Zhang, Y., Bocquet, M., Mallet, V., Seigneur, C. & Baklanov, A. (2012), ‘Real-time air quality forecasting, part i: History, techniques, and current status’, *Atmospheric Environment* **60**, 632–655.

Zheng, Y., Yi, X., Li, M., Li, R., Shan, Z., Chang, E. & Li, T. (2015), Forecasting fine-grained air quality based on big data, *in* ‘Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining’, pp. 2267–2276.

534 Zimmerman, N., Presto, A. A., Kumar, S. P., Gu, J., Hauryliuk, A., Robinson, E. S., Robin-
535 son, A. L. & Subramanian, R. (2018), ‘A machine learning calibration model using random
536 forests to improve sensor performance for lower-cost air quality monitoring.’, *Atmospheric*
537 *Measurement Techniques* **11**(1).