

# Analysis and Comparison of Missing Values Imputation Methods for Atmospheric Pollution Data

**Jakub Jasinski**

Faculty of Electrical Engineering Warsaw University of Technology Warsaw, Poland

\*Corresponding author: Jakub Jasinski, Faculty of Electrical Engineering Warsaw University of Technology Warsaw, Poland.

Submitted: 19 September 2025 Accepted: 25 September 2025 Published: 02 October 2025

 <https://doi.org/10.63620/MKJAVDIM.2025.1001>

**Citation:** Jasinski, J. (2025). Faculty of Electrical Engineering Warsaw University of Technology Warsaw, Poland. *J Aut Veh Dro and Int Mob*, 1(1), 01-05.

## Abstract

Missing values frequently occur in real-world time series datasets, significantly affecting the precision and reliability of data analysis and machine learning models. This research project aims to explore the types of missing data occurrences and examine various imputation methods. The approaches considered will range from simple statistical techniques to more complex methods such as regression models, neural networks, and LSTM models. The effectiveness of these imputation techniques will be assessed using atmospheric pollution data, with a particular focus on PM10 and PM2.5 levels. Each method's performance will be evaluated based on accuracy, consistency, and its impact on subsequent predictive models. The findings indicate that LSTM models are the most effective, while regression and MLP models, though less accurate, offer faster alternatives. Conversely, mean imputation results in the highest error values.

**Keywords:** Missing Data, Data Imputation, Time Series, Regression, LSTM, Atmospheric Pollution, PM10, PM2.5.

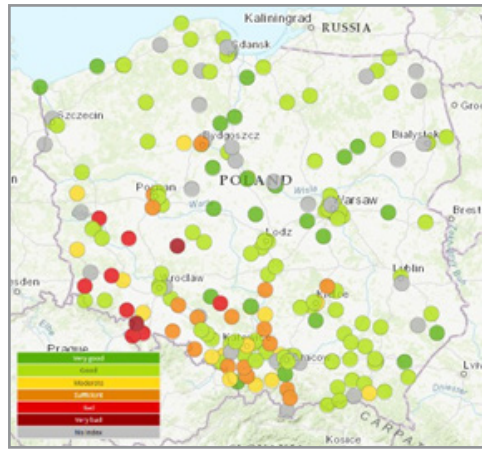
## Introduction

Air pollution is widely recognized as one of the leading threats to both environmental stability and public health today. It stems from various human activities — industrial processes, transportation, agriculture — and natural events like wildfires and dust storms. Regular exposure to polluted air has been shown to cause various health problems, especially affecting the lungs and heart. Fine particles, known as PM2.5, are particularly dangerous because they are small enough to enter deep into the lungs and even reach the bloodstream. This can trigger inflammation throughout the body and worsen conditions such as asthma, chronic bronchitis, or cardiovascular diseases [1]. In recent years, research has also suggested a link between air pollution and damage to the brain, potentially contributing to cognitive decline and developmental disorders in children [2].

Beyond its health effects, air pollution also harms the environ-

ment. It plays a role in acid rain, reduces crop yields, and accelerates climate change. For these reasons, gaining a deeper understanding of how pollution works and how it affects us is key to creating solutions that protect people and the planet alike.

The Polish air quality index is one of the most important indicators for determining the level of air quality in Poland. It is calculated based on 1-hour results of measurements of the following air concentrations: sulfur dioxide (SO<sub>2</sub>), nitrogen dioxide (NO<sub>2</sub>), particulate matter (PM10), particulate matter (PM2.5), and ozone (O<sub>3</sub>). The map presented in Figure 1 displays air quality levels across Poland using a colour-coded index ranging from green (very good) to dark red (very bad). The highest concentrations of poor air quality (red and orange markers) are observed in the southern and southwestern regions, particularly near Wrocław and the Czech border.



**Figure 1: Polish Air Quality Index**

Imputing missing values in air pollution data is essential to ensure accurate monitoring of environmental conditions and detection of harmful pollution levels. Incomplete data can lead to incorrect health risk assessments and hinder timely public health responses. Predictive models for air quality and climate heavily depend on consistent data, and missing values can reduce their accuracy and reliability. Additionally, complete datasets support better policy-making and more effective strategies for pollution control and public safety.

## Literature Review

### Missing Values in Time Series

Missing values in time series refer to absent or unrecorded observations at specific time points in a temporal dataset, which can disrupt the continuity and integrity of time-dependent analysis. These gaps may arise due to sensor failures, transmission errors, or irregular data collection and pose significant challenges for forecasting, modelling, and anomaly detection [3].

### Type of Missing Data

We can distinguish three types of missing data.

- **Missing Completely at Random (MCAR):** Missing values occur independently of both observed and unobserved data. In this case, the probability of data being missing is the same for all observations.

**Example:** A random sensor logging error causes PM10 values to be missed every 1000th reading, regardless of environmental conditions or pollution levels.

- **Missing at Random (MAR):** The missingness is related only to observed data. This means other variables in the dataset can explain the probability of missing values.

**Example:** PM2.5 measurements are frequently missing during rainy weather, which is logged in the dataset. Thus, the missingness can be modelled using weather conditions.

- **Missing Not at Random (MNAR):** The probability of missingness depends on the unobserved data, making it more challenging to model or impute.

**Example:** NO2 concentrations become extremely high during industrial incidents, and the sensors tend to malfunction precisely at these critical levels. The missing data depends on the unmeasured extreme values themselves.

### Dealing With Missing Values

Handling missing values is a critical step in data preprocessing, as it can significantly influence the outcome of any analysis or

model [4]. Below

- **Ignoring:** One of the simplest approaches is ignoring missing values during analysis. This method involves proceeding with computations without accounting for the missing entries. However, we must be sure that omitting these values will not cause the models to malfunction.
- **Deletion:** This method involves removing either the variables (columns) or observations (rows) that contain missing values. While it ensures a clean dataset, deletion can result in significant information loss and reduced statistical power, particularly if missingness is widespread.
- **Imputation:** It refers to filling in missing values with estimated ones based on available data. This approach retains the dataset's structure and size. It can range from basic statistical methods to advanced techniques involving regression models, machine learning, or deep learning algorithms.

## Methodology

### Mean Imputation

In a mean substitution, the mean value of a variable is used in place of the missing data value for that same variable. The theoretical background of the mean substitution is that the mean is a reasonable estimate for a randomly selected observation from a normal distribution. However, with missing values that are not strictly random, especially in the presence of great inequality in the number of missing values for the different variables, the mean substitution method may lead to inconsistent bias [5].

### Regression

Regression imputation using an iterative imputer with Bayesian Ridge regression estimates missing values by modelling each variable as a linear function of the others in multiple rounds. Bayesian Ridge adds regularization through priors, improving stability and handling multicollinearity. This iterative approach refines estimates with each pass, but it may still introduce bias or underestimate variability, especially if data relationships are nonlinear or the model assumptions are violated [6].

### Neural Networks

Using a Multi-Layer Perceptron (MLP) for data imputation involves framing the problem as a supervised learning task, where the MLP learns to predict missing values based on patterns in available historical data. In time series, a common approach is to use a sliding window to transform sequences of past observations into input features for the model [7]. The MLP, being

a feedforward neural network, captures complex nonlinear relationships in the data but does not inherently model temporal dependencies, which can be partially addressed through input engineering. Although MLPs are simpler than recurrent models like LSTM, they are efficient and effective for datasets with short-term dependencies and relatively low missingness.

### LSTM Recurrent Deep Network

LSTM networks introduce a memory cell and a set of gating mechanisms—input, output, and forget gates—that regulate the flow of information. These gates allow the network to retain or discard information over time, making LSTM models particularly effective for tasks involving time-series forecasting, natural language processing, and signal classification [8]. LSTM networks are widely used for imputing missing values in time series by transforming the data into lagged input-output pairs, enabling the model to predict missing values based on temporal context. They are particularly effective in capturing both short- and long-term dependencies and outperform traditional methods under moderate to high missingness. However, LSTMs require careful tuning and more computational resources and struggle

with input sequences that contain missing values unless preprocessed [9].

### Dataset

The GIOS Air Quality Archive provides access to a comprehensive database of air pollution measurements collected across Poland by the State Environmental Monitoring network. This includes two key pollutants of health and environmental concern:

- PM10 (Particulate Matter  $\leq 10$  micrometres): Coarse particles that can penetrate the respiratory system, contributing to respiratory issues and cardiovascular diseases.
- PM2.5 (Particulate Matter  $\leq 2.5$  micrometres): Fine particles capable of reaching the lungs and bloodstream, strongly linked to heart and lung conditions.

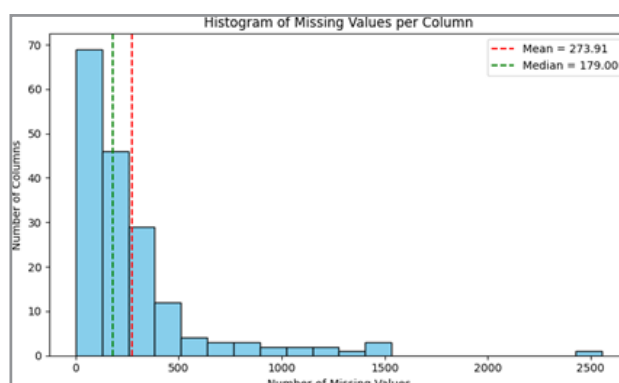
These data are included as time series measured every hour at about 170 stations (depending on the year) in 2000-2023. An analysis was carried out to realize how much data is missing in the data below, resulting in the histograms presented in Figures 2 and 3.

**Table 1:** Comparison of Maximum Missing Interval Statistics (PM10 & PM2.5)

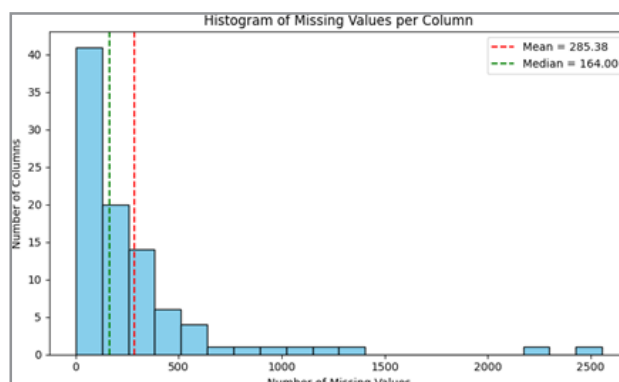
Pollution	PM <sub>10</sub>	PM <sub>2.5</sub>
Max Interval (hrs)	2426.00	2426.00
Min Interval (hrs)	1.00	2.00
Mean Interval (hrs)	140.36	189.36
Median Interval (hrs)	64.00	74.50
Count of Stations	176.00	92.00

As shown in Table I both datasets show significant missing intervals, with some gaps lasting over 100 days, which can severely impact imputation accuracy [10]. Long and irregular gaps

challenge standard methods like mean imputation or iterative imputation significantly when multiple stations are affected simultaneously.



**Figure 2:** Missing Values Count in PM10 Data of Year 2023



**Figure 3:** Missing Values Count in PM2.5 Data of Year 2023

As seen in both cases, most of the data contains less than 250 missing data. The average number of missing data is about 280 in both pollution cases, which is about 3% of the year. However, there are also cases of stations where the number of missing data exceeds 1460, equivalent to 2 months of no measurements and may significantly impact subsequent data analyses. Another critical aspect of missing data analysis is the time intervals in which data is missing. It is much easier to predict a single missing value than a longer interval because of the possibility of taking information from the context of the surrounding measurements.

## Results

To check the accuracy of the previously mentioned methods,

randomly selected values from the available time series will be marked as artificially missing when using the imputation method. Then, the new values will be compared with metrics: SMAPE, MAE, and RMSE [11].

Tables II–IV present the performance of four imputation methods: Mean, Regression, MLP, and LSTM, applied to PM10 and PM2.5 datasets under varying levels of artificially induced missingness: 0.5%, 3%, and 20%. These levels correspond to typical missing durations in air quality datasets, with 0.5% (2 days) representing near-median gap lengths, 3% (11 days) approximating the average number of missing observations, and 20% (2 months) capturing the most extended missing intervals.

**Table 2:** Performance Metrics for PM10 and PM2.5 and 0.5% Missing Data

Model	PM <sub>10</sub>			PM <sub>2.5</sub>		
	SMAPE (%)	MAE (µg/m <sup>3</sup> )	RMSE (µg/m <sup>3</sup> )	SMAPE (%)	MAE (µg/m <sup>3</sup> )	RMSE (µg/m <sup>3</sup> )
Mean	31.93	5.52	9.12	32.10	4.50	7.24
Reg	20.81	3.05	4.78	15.77	2.37	5.63
MLP	19.07	2.78	3.75	16.31	2.15	4.20
LSTM	9.88	2.09	2.83	9.06	1.95	2.60

At minimal missingness, the LSTM model yields the most accurate predictions across all metrics for both pollutants. This demonstrates the model's capability to effectively utilize recent

temporal context for short-gap recovery, outperforming statistical approaches.

**Table 3:** Performance Metrics for PM10 and PM2.5 and 3% Missing Data

Model	PM <sub>10</sub>			PM <sub>2.5</sub>		
	SMAPE (%)	MAE (µg/m <sup>3</sup> )	RMSE (µg/m <sup>3</sup> )	SMAPE (%)	MAE (µg/m <sup>3</sup> )	RMSE (µg/m <sup>3</sup> )
Mean	32.24	5.56	9.14	34.77	5.67	9.74
Reg	21.35	3.33	5.82	15.37	2.29	4.79
MLP	21.57	3.10	4.68	15.62	2.50	4.83
LSTM	11.76	2.56	3.38	11.34	2.42	3.25

When the missing rate increases to a level of dataset average loss, LSTM continues to show superior accuracy, particularly in SMAPE and RMSE. At the same time, MLP also performs

competitively, highlighting the benefits of non-linear learning in moderate data loss scenarios.

**Table 4:** Performance Metrics for PM10 and PM2.5 with 20% Missing Data

Model	PM <sub>10</sub>			PM <sub>2.5</sub>		
	SMAPE (%)	MAE (µg/m <sup>3</sup> )	RMSE (µg/m <sup>3</sup> )	SMAPE (%)	MAE (µg/m <sup>3</sup> )	RMSE (µg/m <sup>3</sup> )
Mean	33.31	5.65	9.41	36.28	5.86	10.72
Reg	21.98	3.36	5.67	16.06	2.35	4.90
MLP	21.76	3.47	5.96	17.05	2.63	5.31
LSTM	14.37	3.11	4.18	14.12	3.10	3.10

With a missing rate representative of extreme data gaps, all methods experience reduced accuracy, yet LSTM maintains its lead across all metrics. This confirms its effectiveness in learning long-term temporal dependencies, making it the most robust method under severe missingness conditions.

## Conclusions

The article presents and compares different solutions of the system forecasting missing data of air pollution PM10 and PM2.5.

An essential point of this research is the analysis of the occurrence of missing data, their length and quantity, and testing methods on different configurations of missing data.

The experimental results reveal that LSTM models consistently deliver the best performance, particularly in moderate to severe missingness scenarios. LSTM's capacity to learn and leverage temporal dependencies allows it to achieve the lowest errors across nearly all conditions, validating its suitability for time

series imputation tasks in environmental datasets.

In contrast, mean imputation performs poorly across all levels of missingness, exhibiting the highest errors and failing to capture even basic temporal structure. This result highlights the risks of relying on naive statistical methods in contexts where the data exhibits strong time-dependent behaviour.

Regression-based and MLP imputations, while not as effective as neural models, consistently outperform mean imputations and show particular promise in low and moderate missingness scenarios. It represents a viable lightweight alternative when computational resources are limited, or model interpretability is prioritized.

Additionally, during the testing process, analyzing individual errors and imputation processes for individual research stations, it can be noticed that using the same method on data from different stations, the error results can differ significantly. To carry out further research, one can focus on more complex relationships in terms of patterns of missing data occurrence, which can help capture the models' ability to apply more effective imputation.

## References

1. Pope, C. A., & Dockery, D. W. (2006). Health effects of fine particulate air pollution: Lines that connect. *Journal of the Air & Waste Management Association*, 56(6), 709–742. <https://doi.org/10.1080/10473289.2006.10464485>.
2. Block, M. L., & Calderón-Garcidueñas, L. (2009). Air pollution: Mechanisms of neuroinflammation and CNS disease. *Trends in Neurosciences*, 32(9), 506–516. <https://doi.org/10.1016/j.tins.2009.05.009>.
3. Pratama, I., Permanasari, A. E., Ardiyanto, I., & Indrayani, R. (2020). A review of missing values handling methods on time-series data. In 2020 3rd International Conference on Information and Communications Technology (ICOIACT) 286–291. IEEE. <https://doi.org/10.1109/ICOIACT50329.2020.9332008>.
4. Salgado, C. M., Azevedo, C., Proença, H., & Vieira, S. M. (2016). Missing data. In *Secondary analysis of electronic health records*, 143–162. Springer.
5. Kang, H. (2013). The prevention and handling of the missing data. *Korean Journal of Anesthesiology*, 64(5), 402–406. <https://doi.org/10.4097/kjae.2013.64.5.402>.
6. Little, R. J., & Rubin, D. B. (2019). *Statistical analysis with missing data* (3rd ed.). Wiley. ISBN: 978-1-118-62262-2. <https://doi.org/10.1002/9781119482260>.
7. Jerez, J. M., Molina, I., García-Laencina, P. J., Alba, E., Ribelles, N., Martín, M., & Franco, L. (2010). Missing data imputation using statistical and machine learning methods in a real breast cancer problem. *Artificial intelligence in medicine*, 50(2), 105–115. <https://doi.org/10.1016/j.artmed.2010.05.002>.
8. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>.
9. Che, Z., Purushotham, S., Cho, K., Sontag, D., & Liu, Y. (2018). Recurrent neural networks for multivariate time series with missing values. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI)*. AAAI Press, 689–696.
10. Chief Inspectorate of Environmental Protection. (2025). Air quality measurement data archive. <https://powietrze.gios.gov.pl/pjp/archives>.
11. Botchkarev, A. (2019). A new typology design of performance metrics to measure errors in machine learning regression algorithms. *Interdisciplinary Journal of Information, Knowledge, and Management*, 14, 45–76, ISSN: 1555-1237. <https://doi.org/10.28945/4184>.