

1 *Conference Proceedings Paper*

2 **Application of a machine learning methodology for** 3 **data implementation**

4 **Chris G. Tzanis ***, Anastasios Alimissis and Ioannis Koutsogiannis

5 Climate and Climatic Change Group, Section of Environmental Physics and Meteorology, Department of
6 Physics, National and Kapodistrian University of Athens, Greece
7 chtzanis@phys.uoa.gr (C.G.T.), alimiss@phys.uoa.gr (A.A.), koutsog@phys.uoa.gr (I.K.)

8 * Correspondence: chtzanis@phys.uoa.gr

9 Received: date; Accepted: date; Published: date

10 **Abstract:** An important aspect in environmental sciences is the study of air quality, using statistical
11 methods (environmental statistics) which utilize large datasets of climatic parameters. The air
12 quality monitoring networks that operate in urban areas provide data of the most important
13 pollutants, which via environmental statistics can be used for the development of continuous
14 surfaces of pollutants concentrations. Generating ambient air quality maps can help guide policy
15 makers and researchers to formulate measures to minimize the adverse effects. The information
16 needed for a mapping application can be obtained by employing spatial interpolation methods to
17 the available data, for generating estimations of air quality distributions. This study uses point
18 monitoring data from the network of stations that operates in Athens. A machine learning scheme
19 will be applied as a method to spatially estimate pollutants' concentrations and the results can be
20 effectively used to implement missing values and provide representative data for statistical analyses
21 purposes.

22 **Keywords:** artificial neural networks; shallow neural networks; machine learning; spatial
23 interpolation; data implementation; air quality

25 **1. Introduction**

26 Studying the distribution of air quality parameters is an important task of urban communities.
27 According to the European Environmental Agency (EEA), air pollution is identified as a major
28 environmental health hazard in Europe as hundreds of thousands of Europeans are affected each
29 year by air quality issues [1-3]. Effective planning strategies require constant monitoring of the
30 various pollutants, creating databases suitable for statistical analysis. Increased data availability can
31 help researchers produce more reliable results. Spatial interpolation techniques have been widely
32 used in air quality studies [4-5] as they can be utilized for data implementation in pollutant time
33 series with missing values and even for sites of interest with no data availability. Additionally, by
34 using these implemented databases, the development of informational tools such as Air Quality
35 Indices (AQI) can be beneficial for presenting in a comprehensible manner new insight to policy
36 makers and the public [6-8]. The EEA proposed a European Air Quality Index (EAQI) which is based
37 on hourly concentrations of five key pollutants (PM₁₀, PM_{2.5}, NO₂, O₃ and SO₂) and has six different
38 levels based on each pollutant's concentrations. This study aims to present a methodology for filling
39 gaps in environmental sciences and specifically in the field of air quality. From the original datasets
40 and based on concentration time series for the selected pollutants of the EAQI, a machine learning
41 data implementation process was followed. This methodology can be utilized as a fast and effective

42 tool which will contribute to the development of indexes such as the EAQI which will subsequently
43 visualize air pollutants' profiles and provide insight in patterns and relationships.

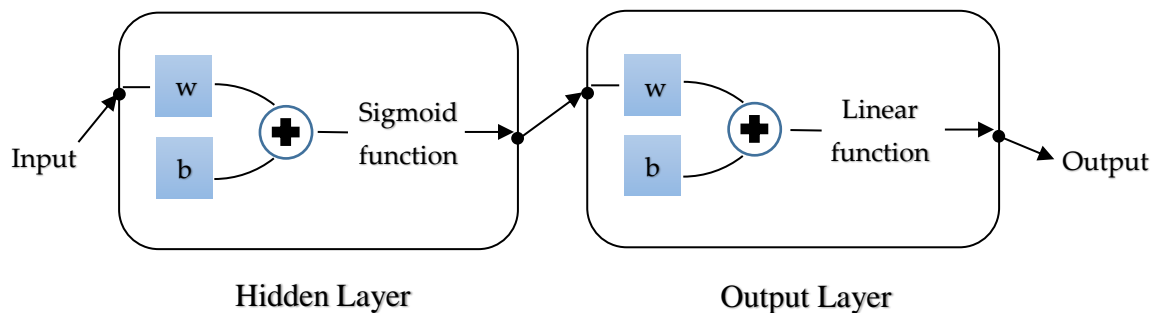
44 2. Experiments

45 2.1 Data

46 The air quality monitoring sites, from which the data were derived, are located at the
47 metropolitan city of Athens, in Greece. As part of the southeastern Mediterranean region, Athens
48 climate is defined by dry summers (long periods, during which the temperatures are considerably
49 high) and wet winters (these periods are usually short) [9]. The basin is bounded by mounts Parnitha,
50 Pentelikon, Hymmetus and Aigaleo to the north, northeast, east-central and west respectively. Due
51 to the transport mechanisms, the topography of the area and the proximity to the sea, the air pollution
52 fields are greatly affected by various flows of different scales [10-13]. The monitoring sites in the area,
53 are part of an air quality monitoring network that operates since 1984, under supervision of the
54 Hellenic Ministry of Environment and Energy (MEE). The network is considered representative of
55 the pollutants' spatial variability and thus suitable for the application of advanced statistical
56 methodologies. For the development of the EAQI, a different number of stations was selected for each
57 pollutant. The criterion for this selection was that a station should have at least a small percent of
58 available data and thus, could contribute to the data implementation methodology. For the five
59 pollutants, NO₂, O₃, PM₁₀, PM_{2.5} and SO₂, the number of stations used was fourteen, thirteen, eleven,
60 six and six, respectively. All five were monitored hourly, and the time period of the analysis was
61 three years (2016-2018).

62 2.2 Methodology

63 The first step in this study, after the database development, was to find the number of gaps that
64 are present in each station's data (target station/missing hourly concentrations) for 2018. This task
65 was performed for all pollutants individually. However, in order to be able to apply effectively the
66 machine learning spatial interpolation scheme, a specific criterion was adopted. For each one of these
67 gaps at a target station, at the same time all the remaining stations must have an available
68 measurement. Even if one of them had also a gap, it was not included in the interpolation process.
69 The results of this step are presented in Table 1 and reveal the number of missing values that can be
70 potentially estimated and used to increase the available data points. The next step was to apply an
71 Artificial Neural Network (ANN) approach for spatial estimation purposes. To achieve this, a
72 Shallow Neural Network (SNN) was utilized as a practical and fairly simple ANN that is moderately
73 demanding in terms of time and computational power. However, it can effectively simulate complex
74 nonlinear relationships between parameters. In detail, two-layer networks with sigmoid hidden
75 neurons and linear output neurons were used (Figure 1).



76 **Figure 1.** A two-layer network with sigmoid hidden neurons and linear output neurons

77 The number of hourly concentrations that were used for the models were those for which none
 78 of the stations had a missing value. The training of the networks was performed with the Levenberg-
 79 Marquardt backpropagation algorithm. The dataset was divided into three subsets used for training,
 80 validation and testing randomly and each subset corresponded to specific percentages of the original
 81 data (70% training, 15% validation, 15% testing). Depending on the pollutant, the number of data
 82 points used for the subsets was different and is presented in Table 2. The network architecture
 83 includes a number of inputs equal to the number of all stations minus the target station (13 for NO₂,
 84 12 for O₃, 10 for PM₁₀, 5 for PM_{2.5} and 5 for SO₂), while the output is always one (target station).
 85 Regarding the number of neurons in the hidden layer, the performance of each network was
 86 evaluated by using the Mean Absolute Error (MAE) statistical criterion [14-18], which is calculated
 87 by using the following equation:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |E_i - O_i| \quad (1)$$

88 where E denotes the estimated concentration, O the observed concentration and n the number of data
 89 points. Lower MAE values illustrate the optimum performing network. Five runs were performed
 90 for all schemes and for hidden layer neurons that ranged from 1 to 40. The best performing networks
 91 and their architecture are presented in Table 3. By using these selected SNN models for the
 92 corresponding inputs of 2018, the gaps in each station and pollutant were filled.

93 3. Results

94 A total of 12,526 missing values were estimated and the percentage of gaps that were filled out
 95 in each station was above 40% for PM₁₀ and PM_{2.5}, above 20% for O₃ and SO₂ and above 15% for NO₂.
 96 Regarding O₃ and NO₂ where the percentage of interpolated values is lower, it needs to be considered
 97 that they had a higher number of input stations and thus, the criterion that none of the inputs should
 98 have a missing value for each gap of the target station, was more difficult to fulfill. Table 1 presents
 99 in detail the gaps originally and after the interpolation, as well as the percentage of missing values
 100 that were estimated.

101

102 **Table 1.** Number of missing values (gaps) during 2018, for the original and spatially interpolated
 103 dataset

	Original gaps	Gaps after interpolation	Difference	Estimated percentage (%)
NO ₂	13,253	11,145	2,108	15.91
O ₃	10,814	7,961	2,853	26.38
PM ₁₀	7,182	3,948	3,234	45.03
PM _{2.5}	4,558	2,524	2,034	44.62
SO ₂	7,043	4,746	2,297	32.61

104

105 The number of data points for the training, validation and testing subsets and for each pollutant
 106 are presented in Table 2. Pollutants with lower number of input stations are associated with higher
 107 data points numbers per station (smaller probability for all the stations to have a missing value at
 108 the same time). However, more stations (NO₂, O₃) provide additional data points. NO₂ and PM_{2.5} are the
 109 pollutants which provided more data for training, validation and testing purposes.

110 **Table 2.** Number of data points distributed to the training, validation and testing subset for the 2016-
 111 2017 time period

	Training	Validation	Testing	Total
NO ₂	47,151	10,101	10,101	67,353
O ₃	25,272	5,412	5,412	36,096
PM ₁₀	13,410	2,880	2,880	19,170
PM _{2.5}	37,785	8,100	8,100	53,985
SO ₂	13,925	3,080	3,080	20,085

112 The architecture of the optimum performance models is presented in Table 3. The hidden
 113 neurons number is an average of all the stations for each pollutant. The MAE average values
 114 (measured in the same units as the concentrations of the pollutants, µg/m³) in these cases are also
 115 included. However, all pollutant-specific networks have the same number of inputs and all networks
 116 have a single output (the target station). The average hidden neuron value ranges from 21.7 to 25.2
 117 which reveals that the models are at an almost equal complexity level.

118 **Table 3.** Number of input, hidden (average) and output neurons as well as MAE (average), mean
 119 concentration values and percentage of error (MAE to mean concentration) for the best performing
 120 models and the 2016-2017 time period

	Input neurons	Hidden	Output	MAE	Mean	Error (%)
NO ₂	13	21.7	1	5.80	32.70	17.74
O ₃	12	22.3	1	6.86	58.86	11.65
PM ₁₀	10	23.6	1	5.71	29.53	19.34
PM _{2.5}	5	25.2	1	5.17	23.81	21.71
SO ₂	5	22.5	1	1.89	6.06	31.19

121 4. Discussion

122 According to Table 3 results, it can be concluded that the error percentage is higher when the
 123 number of input stations is lower and subsequently the information provided for training is more
 124 limited. O₃ is an exception to this statement because although the number of input stations is 12
 125 versus 13 for NO₂ and correspondingly the available data points are nearly half, the error percentage
 126 is considerably lower. This can be explained by examining other behavioral characteristics of this
 127 pollutant (differences in mean values among stations, more easily identifiable patterns in datasets
 128 etc.). When comparing PM_{2.5} and SO₂, where the input neurons are five for both, the prediction
 129 performance for SO₂ is lower, possibly due to the smaller number of data points, according to Table
 130 2 (PM_{2.5} has nearly three times more data points). Different approaches to evaluate the performance
 131 of the models can be followed (scatter diagrams, correlation metrics, etc.) as well as more types of
 132 similar complexity neural network models can be examined.

133 5. Conclusions

134 This study applied SNN models as a tool for point spatial interpolation of air quality parameters,
 135 using data from an air quality monitoring network located at a densely populated urban area. Five
 136 air quality parameters were selected (PM₁₀, PM_{2.5}, NO₂, O₃ and SO₂), due to their importance in the
 137 field of air quality indexes and more specifically, based on the EAQI (proposed by EEA). The results
 138 highlight that the models' performance is significantly affected by the density of the air quality
 139 monitoring network (number of stations and data points per station) as well as the specific patterns
 140 that characterize each pollutant's concentrations. The training dataset is crucial for the networks'
 141 development and needs to be carefully selected in order to provide adequate information which will
 142 augment the networks' generalization ability. This work can be utilized as an alternative for

143 commonly used spatial interpolation methods in the field of air quality and further improvements
144 can be made by using more advanced networks and/or adding meteorological parameters as inputs.

145 **Author Contributions:** C.G.T. and A.A. were involved into the conceptualization, writing-original draft
146 preparation and writing-review and editing of this work, while individually C.G.T. was responsible for the data
147 curation, validation of the results and supervised the whole procedure. All authors performed the various steps
148 of the methodology, processed the data and developed the neural network models. All authors were involved
149 in the discussion of the results and commented on the manuscript. All authors have read and agreed to the
150 published version of the manuscript.

151 **Funding:** This research received no external funding.

152 **Acknowledgments:** The authors would like to acknowledge the Ministry of Environment and Energy for
153 providing the air quality parameters' database which was utilized in this study.

154 **Conflicts of Interest:** The authors declare no conflict of interest.

155 References

- 156 1. Amanollahi, J., Tzani, C., Abdullah, A. M., Ramli, M. F. & Pirasteh, S. (2013). Development of the models
157 to estimate particulate matter from thermal infrared band of Landsat Enhanced Thematic Mapper.
158 *International Journal of Environmental Science and Technology*, 10 (6), 1245–1254.
- 159 2. Baklanov, A., Molina, L. T., & Gauss, M. (2016). Megacities, air quality and climate. *Atmospheric*
160 *Environment*, 126, 235–249. <https://doi.org/10.1016/j.atmosenv.2015.11.059>.
- 161 3. European Environment Agency. (2013). Air quality in Europe—2013 Report: EEA report no 9/2013. In
162 European Union. Retrieved from <http://www.eea.europa.eu/publications/air-quality-in-europe-2013>.
- 163 4. Can, A., Dekoninck, L., & Botteldooren, D. (2014). Measurement network for urban noise assessment:
164 Comparison of mobile measurements and spatial interpolation approaches. *Applied Acoustics*, 83, 32–39.
165 <https://doi.org/10.1016/j.apacoust.2014.03.012>.
- 166 5. Denby, B., Sundvor, I., Cassiani, M., de Smet, P., de Leeuw, F., & Horálek, J. (2010). Spatial Mapping of
167 Ozone and SO₂ Trends in Europe. *Science of the Total Environment*, 408(20), 4795–4806.
168 <https://doi.org/10.1016/j.scitotenv.2010.06.021>.
- 169 6. Zhan, D., Kwan, M. P., Zhang, W., Yu, X., Meng, B., & Liu, Q. (2018). The driving factors of air quality index
170 in China. *Journal of Cleaner Production*, 197, 1342–1351. <https://doi.org/10.1016/j.jclepro.2018.06.108>.
- 171 7. Silva, L. T., & Mendes, J. F. G. (2012). City Noise-Air: An environmental quality index for cities. *Sustainable*
172 *Cities and Society*, 4(1), 1–11. <https://doi.org/10.1016/j.scs.2012.03.001>.
- 173 8. Ganguly, N. D., Tzani, C. G., Philippopoulos, K. & Deligiorgi, D. (2019). Analysis of a severe air pollution
174 episode in India during Diwali festival – a nationwide approach. *Atmosfera*, 32(3), 225–236.
- 175 9. Tzani, C. G., Koutsogiannis, I., Philippopoulos, K. & Deligiorgi, D. (2019). Recent climate trends over
176 Greece. *Atmospheric Research*, 230, 104623, [doi:10.1016/j.atmosres.2019.104623](https://doi.org/10.1016/j.atmosres.2019.104623).
- 177 10. Tzani, C. G., Alimissis, A., Philippopoulos, K., & Deligiorgi, D. (2019). Applying linear and nonlinear
178 models for the estimation of particulate matter variability. *Environmental Pollution*, 246, 89–98.
179 <https://doi.org/10.1016/j.envpol.2018.11.080>.
- 180 11. Deligiorgi, D., Philippopoulos, K., Thanou, L., & Karvounis, G. (2009). A Comparative Study of Three
181 Spatial Interpolation Methodologies for the Analysis of Air Pollution Concentrations in Athens, Greece.
182 *AIP Conference Proceedings*, 1203, 445–450.
- 183 12. Tzani, C. & Varotsos, C. A. (2008). Tropospheric aerosol forcing of climate: a case study for the greater
184 area of Greece. *International Journal of Remote Sensing*, 29 (9), 2507–2517.
- 185 13. Varotsos, C., Christodoulakis, J., Tzani, C. & Cracknell, A. P. (2014). Signature of tropospheric ozone and
186 nitrogen dioxide from space: A case study for Athens, Greece. *Atmospheric Environment*, 89, 721–730.
- 187 14. Cort, J. W., & Kenji, M. (2005). Advantages of the mean absolute error (MAE) over the root mean square
188 error (RMSE) in assessing average model performance. *Climate Research*, 30(1), 79–82.
189 <https://doi.org/10.3354/cr00799>.
- 190 15. Alimissis, A., Philippopoulos, K., Tzani, C. G., & Deligiorgi, D. (2018). Spatial estimation of urban air
191 pollution with the use of artificial neural network models. *Atmospheric Environment*, 191, 205–213.
192 <https://doi.org/10.1016/j.atmosenv.2018.07.058>.

- 193 16. Fallahi, S., Amanollahi, J., Tzani, C. G., & Ramli, M. F. (2018). Estimating solar radiation using
194 NOAA/AVHRR and ground measurement data. *Atmospheric Research*, 199(September 2017), 93–102.
195 <https://doi.org/10.1016/j.atmosres.2017.09.006>.
196 17. Rahimpour, A., Amanollahi, J. & Tzani, C. G. (2020). Air quality data series estimation based on machine
197 learning approaches for urban environments. *Air Quality, Atmosphere & Health*,
198 [doi:https://doi.org/10.1007/s11869-020-00925-4](https://doi.org/10.1007/s11869-020-00925-4).
199 18. Mirzaei, M., Amanollahi, J. & Tzani, C. G. (2019). Evaluation of linear, nonlinear, and hybrid models for
200 predicting PM_{2.5} based on a GTWR model and MODIS AOD data. *Air Quality, Atmosphere & Health*, 12 (10),
201 1215–1224.



© 2020 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).