

EGU22-1255

<https://doi.org/10.5194/egusphere-egu22-1255>

EGU General Assembly 2022

© Author(s) 2022. This work is distributed under the Creative Commons Attribution 4.0 License.



## A Deep Learning approach to de-bias Air Quality forecasts, using heterogeneous Open Data sources as reference

Antonio Pérez<sup>1</sup>, Mario Santa Cruz<sup>1</sup>, Johannes Flemming<sup>2</sup>, and Miha Razinger<sup>2</sup>

<sup>1</sup>Predictia Intelligent Data Solutions S.L., Santander, Spain (predictia@predictia.es)

<sup>2</sup>European Centre for Medium-Range Weather Forecasts (ECMWF), Reading, United Kingdom

The degradation of air quality is a challenge that policy-makers face all over the world. According to the World Health Organisation, air pollution causes an estimate of 7 million premature deaths every year. In this context, air quality forecasts are crucial tools for decision- and policy-makers, to achieve data-informed decisions.

Global forecasts, such as the Copernicus Atmosphere monitoring service model (CAM5), usually exhibit biases: systematic deviations from observations. Adjusting these biases is typically the first step towards obtaining actionable air quality forecasts. It is especially relevant in health-related decisions, when the metrics of interest depend on specific thresholds.

AQ (Air quality) - Bias correction was a project funded by the ECMWF Summer of Weather Code (ESOWC) 2021 whose aim is to improve CAM5 model forecasts for air quality variables ( $\text{NO}_2$ ,  $\text{O}_3$ ,  $\text{PM}_{2.5}$ ), using as a reference the in-situ observations provided by OpenAQ. The adjustment, based on machine learning methods, was performed over a set of specific interesting locations provided by the ECMWF, for the period June 2019 to March 2021.

The machine learning approach uses three different deep learning based models, and an extra neural network that gathers the output of the three previous models. From the three DL-based models, two of them are independent and follow the same structure built upon the InceptionTime module: they use both meteorological and air quality variables, to exploit the temporal variability and to extract the most meaningful features of the past [ $t-24\text{h}$ ,  $t-23\text{h}$ , ...  $t-1\text{h}$ ] and future [ $t$ ,  $t+1\text{h}$ , ...,  $t+23\text{h}$ ] CAM5 predictions. The third model uses the station static attributes (longitude, latitude and elevation), and a multilayer perceptron interacts with the station attributes. The extracted features from these three models are fed into another multilayer perceptron, to predict the upcoming errors with hourly resolution [ $t$ ,  $t+1\text{h}$ , ...,  $t+23\text{h}$ ]. As a final step, 5 different initializations are considered, assembling them with equal weights to have a more stable regressor.

Previous to the modelisation, CAM5 forecasts of air quality variables were actually biased independently from the location of interest and the variable (on average:  $\text{bias}_{\text{NO}_2} = -22.76$ ,  $\text{bias}_{\text{O}_3} = 44.30$ ,  $\text{bias}_{\text{PM}_{2.5}} = 12.70$ ). In addition, the skill of the model, measured by the Pearson correlation, did not reach 0.5 for any of the variables—with remarkable low values for  $\text{NO}_2$  and  $\text{O}_3$  (on average:  $\text{pearson}_{\text{NO}_2} = 0.10$ ,  $\text{pearson}_{\text{O}_3} = 0.14$ ).

AQ-BiasCorrection modelisation properly corrects these biases. Overall, the number of stations that improve the biases both in train and test sets are: 52 out of 61 (85%) for NO<sub>2</sub>, 62 out of 67 (92%) for O<sub>3</sub>, and 80 out of 102 (78%) for PM<sub>2.5</sub>. Furthermore, the bias improves with declines of -1.1%, -9.7% and -13.9% for NO<sub>2</sub>, O<sub>3</sub> and PM<sub>2.5</sub> respectively. In addition, there is an increase in the model skill measured through the Pearson correlation, reaching values in the range of 100-400% for the overall improvement of the variable skill.