

Comment on “Deep Ensemble Machine Learning Framework for the Estimation of PM_{2.5} Concentrations”

Massimo Stafoggia,¹  Giorgio Cattani,² Carla Ancona,¹ Antonio Gasparri,³ and Andrea Ranzi⁴

¹Department of Epidemiology, Lazio Regional Health Service/ASL Roma 1, Rome, Italy

²Institute for Environmental Protection and Research, Rome, Italy

³Department of Public Health, Environments and Society, London School of Hygiene & Tropical Medicine, London, UK

⁴Environmental Health Reference Centre, Regional Agency for Environmental Prevention of Emilia-Romagna, Modena, Italy

<https://doi.org/10.1289/EHP11385>

Refers to <https://doi.org/10.1289/EHP9752>

Yu et al.¹ applied an innovative methodology, deep ensemble machine learning, to estimate daily concentrations of ambient particulate matter with an aerodynamic diameter of <2.5 μm (PM_{2.5}) in 2015–2019 at 1-km² spatial resolution in Italy. To do so, they trained multiple prediction models on PM_{2.5} concentrations measured at 133 monitoring stations.

We recently published similar studies in Italy using alternative methods, including mixed-effects models,² random forests,^{3,4} ensemble techniques,⁵ and Bayesian approaches.⁶ We believe that Yu et al.¹ did not adequately consider critical questions, negatively affecting the validity and interpretation of their results.

First, the set of monitoring stations measuring fine particles in Italy during 2015–2019 is much larger than the one used by Yu et al. (Figure S2 in their appendix¹). Figure 1 shows the 289 monitors measuring PM_{2.5} during a slightly different period (2016–2019) that used in our recent publication.⁴ The comparison of the two maps shows that Yu et al. did not include stations in Southern Italy or in the two main islands, Sicily and Sardinia. These areas have unique geoclimatic conditions and source profiles of ambient PM_{2.5} concentrations, with a mixture of anthropogenic emissions from large industrial plants and heavily urbanized areas, coupled with natural sources such as sea salt, desert dust from North Africa, forest fires and volcanic emissions from Mount Etna.⁷ Such complexity is extremely difficult to capture with any empirical predictive model.²

Second, the cross-validated coefficient of determination and root mean square errors reported by Yu et al. cannot be compared with those previously published,^{2–6} because they are based on a small number of monitors selected in areas where a higher performance of spatiotemporal prediction models has been previously documented.⁴

Third, we are surprised not to find several key predictors of spatial (e.g., road network, impervious surfaces, industrial sites) or spatiotemporal (e.g., desert dust episodes, outputs from atmospheric dispersion models, planetary boundary layer, vegetation indices) variability of PM_{2.5}. Such predictors were used in previous applications^{2–6} and allowed the capture, at least partially, of the geoclimatic complexity of southern regions.



Figure 1. Map of the 289 PM_{2.5} monitoring sites available in Italy during 2016–2019 PM_{2.5}.

In conclusion, we consider the methodological effort from Yu et al. a valid contribution to the literature. However, because their model represents only 6 of 20 regions contributing sufficient data, we question the use of their PM_{2.5} estimates for later epidemiological studies using Italian data.

References

- Yu W, Li S, Ye T, Xu R, Song J, Guo Y. 2022. Deep ensemble machine learning framework for the estimation of PM_{2.5} concentrations. *Environ Health Perspect* 130(3):37004, PMID: 35254864, <https://doi.org/10.1289/EHP9752>.
- Stafoggia M, Schwartz J, Badaloni C, Bellander T, Alessandrini E, Cattani G, et al. 2017. Estimation of daily PM₁₀ concentrations in Italy (2006–2012) using finely resolved satellite data, land use variables and meteorology. *Environ Int* 99:234–244, PMID: 28017360, <https://doi.org/10.1016/j.envint.2016.11.024>.
- Stafoggia M, Bellander T, Bucci S, Davoli M, de Hoogh K, De'Donato F, et al. 2019. Estimation of daily PM₁₀ and PM_{2.5} concentrations in Italy, 2013–2015, using a spatiotemporal land-use random-forest model. *Environ Int* 124:170–179, PMID: 30654325, <https://doi.org/10.1016/j.envint.2019.01.016>.
- Stafoggia M, Cattani G, Ancona C, Ranzi A. 2020. Exposure assessment of air pollution in Italy 2016–2019 for future studies on air pollution and COVID-19. [In

Address correspondence to Massimo Stafoggia, Department of Epidemiology, Lazio Regional Health Service/ASL Roma 1, Via Cristoforo Colombo, 112-00147 Rome, Italy. Email: m.stafoggia@deplazio.it

The authors declare they have no actual or potential competing financial interests.

Note to readers with disabilities: *EHP* strives to ensure that all journal content is accessible to all readers. However, some figures and Supplemental Material published in *EHP* articles may not conform to 508 standards due to the complexity of the information being presented. If you need assistance accessing journal content, please contact ehpsubmissions@niehs.nih.gov. Our staff will work with you to assess and meet your accessibility needs within 3 working days.

- Italian.] *Epidemiol Prev* 44(5–6 suppl 2):161–168, PMID: 33412807, <https://doi.org/10.19191/EP20.5-6.S2.115>.
5. Shtein A, Kloog I, Schwartz J, Silibello C, Michelozzi P, Gariazzo C, et al. 2020. Estimating daily PM_{2.5} and PM₁₀ over Italy using an ensemble model. *Environ Sci Technol* 54(1):120–128, PMID: 31749355, <https://doi.org/10.1021/acs.est.9b04279>.
 6. Fioravanti G, Martino S, Cameletti M, Cattani G. 2021. Spatio-temporal modelling of daily concentrations in Italy using the SPDE approach. *Atmos Environ* 248:118192, <https://doi.org/10.1016/j.atmosenv.2021.118192>.
 7. Sistema Nazionale per la Protezione dell'Ambiente. 2020. La Qualità dell'Aria in Italia. Edizione 2020. SNPA, Rapporti 17/2020. <https://www.snpambiente.it/wp-content/uploads/2020/12/QUALITA-ARIA-ITALIA.pdf> [accessed 13 April 2022].