

Supporting Information

XIS-PM_{2.5}: A daily spatiotemporal machine-learning model for PM_{2.5} in the contiguous United States

Allan C. Just^{1*}, Kodi B. Arfer¹, Johnathan Rush¹, Alexei Lyapustin², Itai Kloog^{1,3}

¹Department of Environmental Medicine and Public Health, Icahn School of Medicine at Mount Sinai, New York, NY, USA

²NASA Goddard Space Flight Center, Greenbelt, MD, USA

³The Department of Geography and Environmental Development, Ben-Gurion University of the Negev, Beer Sheva, Israel

Corresponding Author: allan.just@mssm.edu

Address: Allan Just, One Gustave L. Levy Place, Box 1057, New York, NY 10029 USA

Table 3: 2019 cross-validation results ($\mu\text{g}/\text{m}^3$) broken down by meteorological season. Results for December are taken from the 2018 model so that a contiguous winter is analyzed.

Season	Observations	Sites	MAD	MAE	Bias
all	349,993	1,284	3.15	1.72	-0.20
Spring	88,936	1,256	2.99	1.55	-0.20
Summer	88,942	1,253	2.81	1.47	-0.13
Fall	87,001	1,245	3.14	1.70	-0.21
Winter	85,114	1,284	3.63	2.16	-0.27

Table 3 shows cross-validation results by season for one year. Compared to the whole one-year period, MAD and MAE are lower in spring and summer and higher in winter.

Table 4: Results from each yearly cross-validation among isolated sites, in $\mu\text{g}/\text{m}^3$.

Year	Observations	Sites	MAD	MAE	Bias
2003	25,913	236	4.67	1.93	-0.03
2004	29,202	231	4.43	2.00	0.07
2005	30,520	225	5.02	2.09	-0.01
2006	34,457	233	4.38	2.08	0.01
2007	37,396	225	4.78	2.23	-0.06
2008	38,129	226	4.24	2.05	-0.20
2009	40,172	223	3.83	1.94	-0.18
2010	44,377	230	3.98	1.98	-0.06
2011	44,935	223	4.18	2.21	-0.25
2012	46,399	222	3.83	2.15	-0.27
2013	49,163	224	3.66	2.10	-0.11
2014	54,287	238	3.59	2.04	-0.26
2015	55,638	242	3.59	2.04	-0.19
2016	57,603	241	3.04	1.86	-0.17
2017	60,450	239	3.49	2.02	-0.13
2018	62,863	239	3.58	1.95	-0.28
2019	63,538	242	3.08	1.73	-0.20
2020	65,667	238	3.49	2.01	-0.22
2021	64,125	242	3.97	2.06	-0.25

In addition to the weighted analyses using all stations, we wished to evaluate performance where ground networks were especially sparse. Thus, Table 4 shows unweighted MAE from cross-validation among the sites that were particularly isolated, defined as being more than 50 km from all other sites available in the same year.

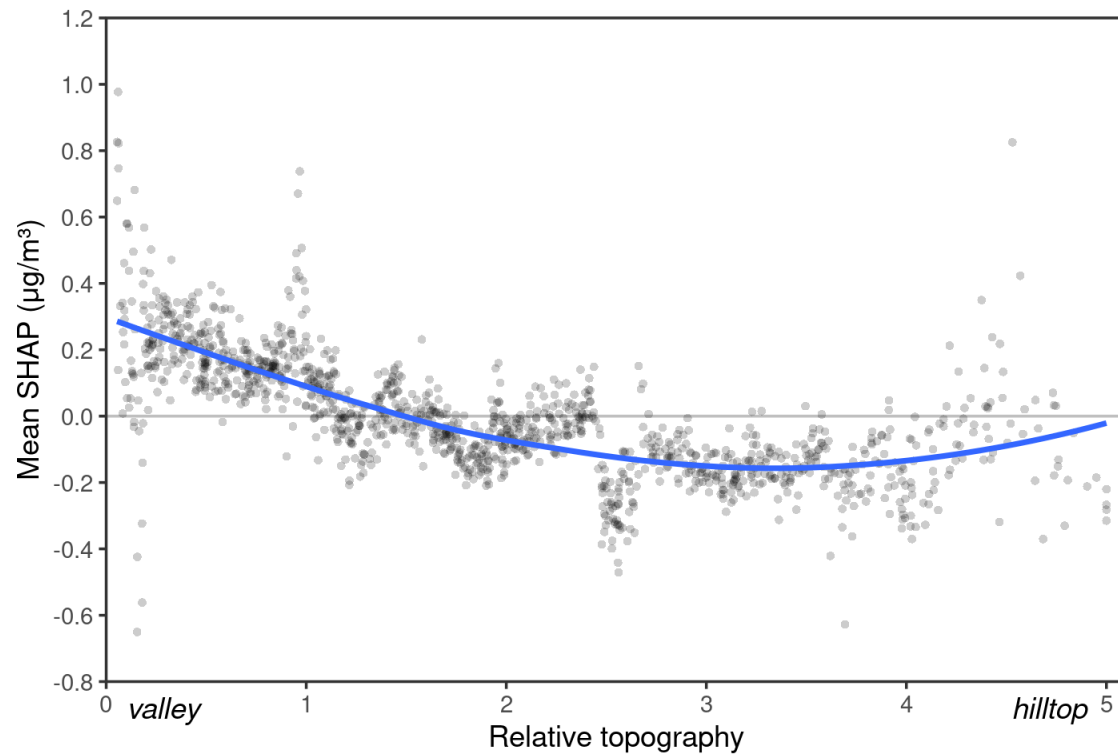


Figure 6: SHAP of hilliness as a function of hilliness.

Figure 6, similar to Figure 3, plots the mean SHAP of the hilliness feature for each site. We see higher average predicted $\text{PM}_{2.5}$ at sites in a valley versus on a hill.

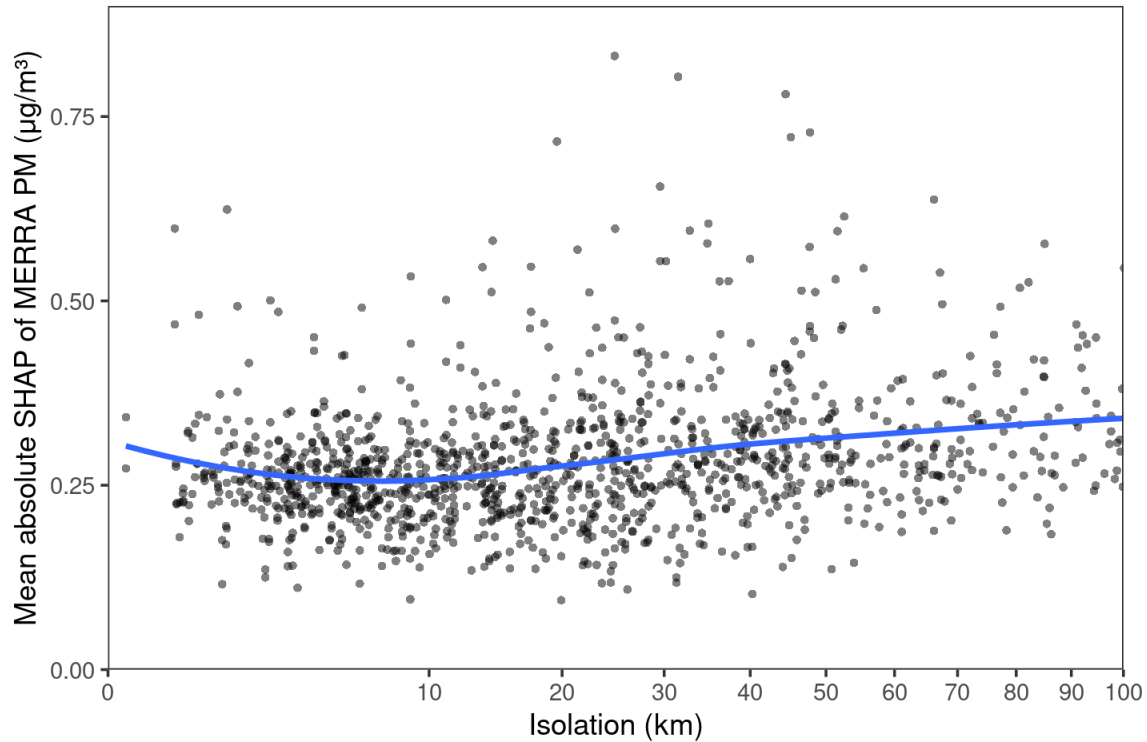


Figure 7: SHAP of modeled surface $\text{PM}_{2.5}$ concentrations from MERRA-2 as a function of site isolation. The x -axis is on a square-root scale.

Figure 7 is an example of the relationship between the SHAP of one variable and a different quantity. The y -axis shows the per-site mean absolute SHAP of MERRA-2 modeled $\text{PM}_{2.5}$, but the x -axis shows the site's distance from its nearest neighbor; that is, its degree of isolation. Similarly, we examined how the SHAP of the IDW feature varied according to isolation. The per-site mean absolute SHAP for IDW in 2010 was Kendall-correlated -0.17 with the distance to the nearest other site, meaning that the IDW is less influential on predictions for more isolated sites.