



Supplement of

Predicting peak daily maximum 8 h ozone and linkages to emissions and meteorology in Southern California using machine learning methods (SoCAB-8HR V1.0)

Ziqi Gao et al.

Correspondence to: Ziqi Gao (zgao71@gatech.edu)

The copyright of individual parts of the supplement might differ from the article licence.

1. Detailed information of emission data

Annual average NOx and VOC emissions from 2000 to 2035 in the South Coast Air Basin (SoCAB) are projected from the emissions in 2012 (Cox et al., 2013). To backcast the NOx and VOC emissions from 1990 to 2000 based on the 2012 emission inventories, first, we computed the NOx and VOC emissions ratios (the emissions were projected from the emissions in 2008) in 1990 and 1995 to the year 2000 (Cox et al., 2013; Cox et al., 2009). Furthermore, we got the adjusted emissions in 1990 and 1995 by multiplying the ratios to the emissions in 2000 that were estimated from the inventory in 2012. Finally, we used linear interpolation to compute the emissions of the years between 1990 and 1995 and between 1995 and 2000.

Kind of Variables	Variables Units		Data Source	
Response Variable	Top 30 MDA8 OzoneConcentrations	ppbV	CARB/ EPA	
	Temperature	°C		
Surface	Wind Speed	m/s	NOAA ¹ / CARB	
Meteorology	Wind Direction	Degree	-	
	Solar Radiation ^b	W/m ²	CARB/ EPA/ NSRD	
	Geopotential Height	m		
	Temperature	°C		
(500 and 850	Dew Point Temperature	°C	NOAA ²	
$(500 \text{ and } 850 - \text{millibar})^{\circ}$	Wind Speed	m/s		
	Wind Direction ^d	Degree		
	Relative Humidity ^e	%	-	
Estimated Emissions ^f	ated NOx/ VOC		CARB	
Large-scale Climate Index	Niño 3.4 monthly indices	°C	CPC	
Temporal Variable	Day of Year	— None	ΝA	
	Day of Week	INUIIC		

Table S1. List of the data sources of the variables used to build the computational models of top 30 MDA8 ozone days from 1990 to 2019.

Data Source Abbreviation: **CARB**: California Air Resources Board (<u>https://www.arb.ca.gov/aqmis2/aqdselect.php</u>, last access May 23, 2020); **EPA**: EPA AQS air pollutant data queries (<u>https://aqs.epa.gov/aqsweb/airdata/download_files.html</u>, last access

May 27, 2020); NOAA¹: National Oceanic and Atmospheric Administration (<u>https://www.arb.ca.gov/aqmis2/metselect.php</u>, last access May 27, 2020); NOAA²: National Oceanic and Atmospheric Administration (<u>https://ruc.noaa.gov/raobs/</u>, last access May 23, 2020); NSRD: National Solar Radiation Database (https://nsrdb.nrel.gov/, last access May 27, 2020); CPC: Climate Prediction Center (https://www.cpc.ncep.noaa.gov/data/indices/sstoi.indices, last access May 23, 2020).

a: All the surface meteorological variables were obtained from Barstow-Daggett Airport and Los Angeles International Airport (LAX).

b: To avoid the outliers in the CARB and EPA dataset and create a continuous solar radiation (SR) from 1990 to 2019, we combined the SR data at LAX from NSRD meteorological statistical model and those at Santa Clarita site/ Los Angeles N Main Street site/ Victorville Park Avenue site from CARB and EPA AQS archives. We implemented the missing SR value using the data at Joshua Tree NP Black Rock site.

c: Upper meteorological data is at the Miramar site, close to the SoCAB, and no site has sounding data in the SoCAB. The upper meteorological data at the Miramar site that follows the standard radiosonde release time is relatively more than other sites (e.g., Edwards Air Force Base (AFB), Vandenberg AFB, Point Mugu, and San Nicolas Island).

d: We used the sine of the wind direction at upper air to represent the transport direction.

e: Relative Humidity (RH) value at 500 and 850 millibar (mb) was computed through the Clausius-Clapeyron Equation (Alduchov et al., 1996; Lawrence, 2005).

$$RH = e^{\{5321 \times \{(\frac{1}{T} - \frac{1}{Td}\}\}}$$
(Equation 1)

where T is air temperature and Td is dew point temperature.

f: The full description of the estimated emissions from 1990 to 2000 is given in the SI: Detailed information of emission data.

Kind of Variables	Variables	Abbreviation	Unit
Temporal	Day of the week (factor, from Mon to Sun)	dayofweek	None
Variables	Day of year (from 1 to 365/366)	dayofyear	None
Surface Meteorological Variables	Daily maximum surface temperature at the Barstow Airport/ LAX site	TmaxBarstow/ TmaxLAX	°C
	Daily minimum surface temperature at the Barstow Airport/ LAX site	TminBarstow/ TminLAX	°C
	Daily average wind speed at the Barstow Airport/ LAX site	AWNDBarstow/ AWNDLAX	m/s
	Max/ Mean solar radiation	SRmax/ SRmean	W/m ²
	Daily RH at 500/ 850 mb	Mir500RH	%

Table S2. Predictors used to test the final/ optimal GAM, MARS, SVR and RF model equations.

		/Mir850RH	
	Daily dew point temperature at	MirDewPtT500C	°C
_	500/ 850 mb	/MirDewPtT850C	C
Upper Air	Daily temperature at 500/850 mb	MirTemp500C /	°C
Meteorological	Daily temperature at 500/ 850 mb	MirTemp850C	
Variables	Daily wind speed at 500/ 850 mb	MirWS500ms/	m/s
_	Daily which speed at 500/ 850 mb	MirWS850ms	111/ S
	Daily wind direction at 500/850	MirWD500/Mir	None
_	mb*	WD850	None
	Daily height at 500/ 850 mb	MirHeight500/	m
	Daily height at 500/ 850 mb	MirHeight850	
Large-scale climate pattern	Monthly Niño 3.4 indices	ENSOmonthly	°C
Emissions —	Annual averaged NOx emissions	eNOx	Tons/day
	Annual averaged VOC emissions	eROG	Tons/day

*: We used the sine of the wind direction at upper air to represent the transport direction.

Table S3. Summary of statistical results of the top 30 MDA8 ozone concentrations using four methods at Crestline site.

Method	Mean Bias (ppbV)	R ²	RMSE (ppbV)	
GAM model	-0.02	0.84	9.74	
MARS model	-0.40	0.83	10.1	
RF model ¹	-0.44	0.81	10.9	
RF model ²	-0.36	0.81	10.9	
SVR model ¹	-1.2	0.81	10.8	
SVR model ¹ +tune	-0.74	0.81	10.9	
SVR model ²	-1.2	0.83	10.4	

1 and 2: RF/ SVR model with the same variables as GAM model and RF/ SVR model with the optimal combination of the indicators.

Table S4. Summary of statistical results of the top 30 MDA8 ozone concentrations using four methods at Crestline site using 10-fold cross validation (90% is training data and 10% is testing data).

Method	Training Data		Testing Data		
	R ²	RMSE (ppbV)	R ² RMSE (ppbV)		
GAM model	0.84	9.74	0.85	9.67	
MARS model	0.83	10.3	0.83	10.2	

RF model	0.80	11.0	0.82	10.3
SVR model ¹	0.81	10.9	0.81	10.4
SVR model ²	0.82	10.4	0.8	10.6

1 and 2: SVR model with the same variables as GAM model and SVR model with the optimal combination of the indicators.

Table S5. Summary of statistical results of the top 15 MDA8 ozone concentrations and the 4^{th} highest ozone predictions using four methods at Crestline site.

	Top 15 MDA8 ozone days		4 th highest MDA8 ozone			
Method [–]	Mean Bias (ppbV)	R ²	RMSE (ppbV)	Mean Bias (ppbV)	R ²	RMSE (ppbV)
GAM	0.02	0.90	8.30	-3.94	0.98	5.64
MARS	-0.27	0.89	8.55	-4.84	0.97	6.76
RF ¹	-0.40	0.85	10.2	-6.09	0.97	8.12
RF ²	-0.24	0.85	10.1	-5.89	0.96	8.39
SVR ¹	-1.22	0.86	9.92	-4.31	0.93	7.37
SVR ²	-1.16	0.88	9.19	-4.60	0.90	9.73



Figure S1. Correlation value among all the available independent variables.



Figure S2. Number of the remaining terms in the built MARS model vs the RMSE value using 10-fold validation with the training dataset (90% of the original dataset). The red point shows the best setting that remain 14 terms in the MARS model and RMSE equals to 10.19 ppbV. The RMSE of the 16 terms MARS model is 10.27 ppbV.



Figure S3. Number of trees in the built RF model vs the RMSE value.



Figure S4. Number of variables in each tree of the built RF model vs the out-of-bag (OOB) value.



Figure S5. Observed (blue) and predicted top 30 MDA8 ozone concentrations using original (orange) and tuned (green) SVR models from 1990 to 2019 at Crestline site.



Figure S6. Observed and predicted top 30 MDA8 ozone concentrations with the corresponding annual NOx and VOC emissions and maximum temperature from 1990 to 2019 at Crestline site (the color of the points shows the value of maximum temperature, annual NOx and VOC emissions).

References

- Alduchov, O. A., & Eskridge, R. E. (1996). Improved Magnus Form Approximation of Saturation Vapor Pressure. *Journal of Applied Meteorology*, 35(4), 601-609.
- Cox, P., Delao, A., & Komorniczak, A. (2013). The California almanac of emissions and air quality.

California Air Resources Board, Sacramento, CA.

- Cox, P., Delao, A., Komorniczak, A., & Weller, R. (2009). The California almanac of emissions and air quality. *California Air Resources Board, Sacramento, CA*.
- Lawrence, M. G. (2005). The relationship between relative humidity and the dewpoint temperature in moist air A simple conversion and applications. *Bulletin of the American Meteorological Society*, 86(2), 225-233.