

Supporting Information

Constrained Mixed-Effect Models with Ensemble Learning for Prediction of Nitrogen Oxide Concentrations at a High Spatiotemporal Resolution

Lianfa Li^{*,1,2}, Fred Lurmann³, Rima Habre^{*,1}, Robert Urman¹, Edward Rappaport¹, Beate Ritz⁴,

Jiu-Chiuan Chen¹, Frank D. Gilliland¹, Jun Wu^{*,5}

1. Department of Preventive Medicine, University of Southern California, Los Angeles, CA, USA

2. State Key Laboratory of Resources and Environmental Information System, Institute of Geographical Sciences and Natural Resources, Chinese Academy of Sciences, Beijing, China

3. Sonoma Technology, Inc., Petaluma, CA, USA

4. Department of Epidemiology, University of California, Los Angeles, CA, USA

5. Program in Public Health, College of Health Sciences, University of California, Irvine, CA, USA

*Corresponding authors' contact information: Lianfa Li: Phone, 650-388-0268; Email, lianfali@usc.edu or lspatial@gmail.com. Rima Habre: Phone, 323-442-8283; Email, habre@usc.edu. Jun Wu: Phone, 949-824-0548; Email, junwu@uci.edu.

Number of pages: 21

Number of figures: 8

Number of tables: 4

Contents

1. Overview	S1
2. Measurements of NO₂ and NO_x concentrations	S1
3. The Covariates Selected	S1
3.1 CALINE4-estimated concentrations from local traffic emissions	S1
3.2 Traffic density	S2
4. Modeling Approach	S2
4.1 Non-parametric additive methods	S2
4.2 Modeling of spatial random effects	S3
4.3 Aggregated predictions by ensemble learning	S3
Table S1. Correlation between average biweekly passive measurements at USC, UCLA, and UCI sites & collocated routine monitoring sites, and their linear regression coefficients used for consistent adjustment	S5
Table S2. Variance explained by the predictors included in the Stage 1 mixed effects model	S6
Table S3. NO ₂ model performance by CHS community in cross validation	S7
Table S4. NO _x model performance by CHS community in cross validation	S8
Figure S1. Study region with routine monitoring locations and USC sampling locations (the UCLA and UCI data are not shown due to IRB restrictions)	S9
Figure S2. Histograms for NO ₂ and NO _x to determine log-transformation (no log transformation for NO ₂ ; log transformation for NO _x).....	S10
Figure S3. The first and second temporal basis functions to reflect the seasonal variability for the study region	S11
Figure S4. Non-linear association between predictive variables and concentrations by mixed models	S12
Figure S5. Spatial topology for spatial effect modeling by Thiessen polygons (aggregate distance: 500 m)	S14
Figure S6. Residual plot between observed values vs. residuals (a. NO ₂ ; b. NO _x) for the sample selected by bootstrap aggregating.....	S15
Figure S7. Spatial distribution of the predicted and observed NO ₂ and NO _x concentration means across the CHS communities	S16
Figure S8. Boxplot for correlation between constrained prediction and observed values for the time series of 51 routine monitoring stations	S17
Figure S9. Time series simulated for the routine monitoring sites with the minimum correlation between constrained prediction and observed values (a.0.55 for NO ₂ ; 0.70 for NO _x).....	S18
Figure S10. Summer (a) and winter (b) averages of the 2005-2006 biweekly NO _x at the CHS subject locations for San Dimas	S20

1 **1. Overview**

2 This paper proposed a three-stage spatiotemporal model that can reliably predict
3 nitrogen oxide concentrations with a high spatiotemporal resolution over a long time
4 span (>20 years). The spatially extensive highly-clustered exposure data from
5 short-term measurement campaigns across 1-2 years and long-term central site
6 monitoring in 1992-2013 were leveraged to develop the first stage mixed-effect
7 models and the second stage ensemble learning with uncertainty estimates. Then at
8 the third stage, constrained optimization was designed and implemented based on the
9 point estimates from the first and second stages to simulate the long-term series of
10 pollutant concentrations for any target location in the study region. The following
11 sections provide the supplemental information about concentration measurements, the
12 covariates selected, modeling method, and results to support the formal paper.

13 **2. Measurements of NO₂ and NO_x concentrations**

14 Besides the measurements of routine monitoring stations, additional data were
15 generated in intensive field measurement campaigns conducted by the University of
16 Southern California (USC), University of California Los Angeles (UCLA), and
17 University of California Irvine (UCI), respectively. Passive diffusion-based Ogawa
18 samplers ¹ were used to measure NO₂ and NO_x at different time periods and at
19 different locations (e.g. outside homes, schools, strategic outdoor sampling locations,
20 and central monitoring sites). Our previous papers respectively provide more details
21 about the measurement methods that generated the USC ², UCLA ³ and UCI ⁴ samples.
22 [Figure S1](#) also shows the locations for the routine and USC sampling sites [the UCI
23 and UCLA sampling locations are concealed to comply with specific requirements by
24 their Institutional Review Boards].

25 To minimize systematic bias in the field campaign data, we compared the passive
26 data with the active data from the routine government monitors at the co-located sites
27 and made small adjustments to standardize the passive measurements to Federal
28 Reference Method equivalent values. For details, please refer to [Table S1](#).

29 **3. The Covariates Selected**

30 **3.1 CALINE4-estimated concentrations from local traffic emissions**

31 CALINE4 is a line source dispersion model that was used to assess the

32 contribution of local motor vehicle emissions to ambient concentrations^{5, 6}.
 33 CALINE4 was used to compute mean NO_x concentration from emissions on freeways
 34 [coded using Feature Class Codes (FCC) as FCC1] and non-freeways (FCC2, FCC3
 35 and FCC4). Traffic count data were obtained from Caltrans and TeleAtlas/GDT and
 36 assigned to ERSI Premium Street Map roadway geometry. Emission strength was
 37 estimated using quarterly average daily traffic volumes and EMFAC2011 (for
 38 1992-2012)⁷ and EMFAC2014 (for 2013)⁸, which generated air basin emission
 39 factors that were based on average vehicle speed and heavy-duty truck fraction
 40 (Caltrans post-mile truck count data by year). Wind speeds and directions were
 41 based on hourly observations of these surface meteorological variables from 72
 42 monitoring stations of California⁸.

43 **3.2 Traffic density**

44 Traffic density represents distance-decayed annual average daily traffic (AADT)
 45 volume in both directions from all roads (FCC1-FCC4) within a circular buffer.
 46 Traffic density is symmetric on both sides of each roadway, computed as if the wind
 47 directions were uniformly distributed around the compass. The values of traffic
 48 density were computed by the ESRI ArcGIS density function using a kernel with a
 49 300 m search radius and 5 m grid resolution. Annual traffic density estimates were
 50 provided for the regulatory monitoring sites and local sampling locations. In
 51 addition, because these cover a long time period, the traffic densities were scaled by
 52 the South Coast Air Basin (SoCAB) EMFAC2011 vehicle fleet average NO_x emission
 53 factor for 50 mph and 6% heavy-duty vehicle fraction (normalized to 1.00 in 2002
 54 since we used the 2002 AADT as the baseline data) to reflect the composite trend in
 55 traffic volumes and emissions over time.

56 **4. Modeling Approach**

57 **4.1 Non-parametric additive methods**

58 For spatiotemporal factors, we adopted non-parametric additive methods to
 59 model non-linear effects. Specifically, we used non-parametric trend functions to
 60 quantify the association, $s(\dots)$, i.e. approximated by the weighted sum of polynomial
 61 spline (B-spline basis) functions

$$62 \quad x_{min} = \zeta_0 < \zeta_1 < \zeta_2 \dots < \zeta_{m-1} < \zeta_m = x_{max} \quad (S1)$$

$$63 \quad f(x_i) = \sum_j \beta_{ij} B_{ij}(x_i) \quad (S2)$$

64 where ζ_i is the split for an interval for the covariate x_i (ζ_0 and ζ_m are respectively the
 65 minimum and maximum values of x_i), B_{ij} and β_{ij} respectively represent the basis
 66 function and parameter for the interval j for the covariate x_i . Penalized maximum

67 likelihood was used to solve (eq. S2) to estimate β_{ij} .

68 4.2 Modeling of spatial random effects

69 By estimating a structured component and an unstructured component, we can
 70 distinguish between the two sources of spatial autocorrelations⁹. The spatial effects
 71 were modeled using the following formulas for both structured (S3) and unstructured
 72 spatial effects (S4).

$$73 \quad f_s(r_s)|f_s(r'), r' \neq r_s, \tau_{str}^2 \sim N\left(\frac{1}{N_{r_s}} \sum_{r' \in \delta_{r_s}} f_s(r'), \frac{\tau_{str}^2}{N_{r_s}}\right) \quad (S3)$$

$$74 \quad f_{re}(r_s)|\tau_{unstr}^2 \sim N(0, \tau_{unstr}^2) \quad (S4)$$

75 where r_s is the region where the observation $y(s,t)$ is located, δ_{r_s} represents a set of
 76 neighbors (r') of the polygon r_s , N_{r_s} is the number of neighboring polygons for r_s ,
 77 τ_{str}^2 is the total variance for the structured component, $\tau_{str}^2 \sim \text{IG}(a,b)$. $f_s(r')$ in (eq.
 78 S3) represents the spatial influence from neighboring polygons (r') on r_s ; $f_{re}(r_s)$ in (eq.
 79 S4) represents the unstructured spatial effect with zero mean and standard deviation
 80 (τ_{unstr}^2) for r_s .

81 Thiessen polygons were constructed around the central points (derived by
 82 averaging the coordinates of all the routine or/and campaigns sampling locations
 83 within a certain distance) to simulate spatial effects. Rook adjacency was used for
 84 spatial adjacency: two polygons were assumed to be neighbors if they share a
 85 common border. We conducted sensitivity tests for a series of aggregation distances
 86 (100 m, 300 m, 500 m and 3 km) and finally selected an optimal aggregate distance
 87 (500 m) that provided a good balance between model accuracy and computing
 88 efficiency. We used the packages of `rgdal` and `spdep` in the statistics software R
 89 (Version 3.3) for generation of the Thiessen polygons with their spatial weight matrix.

90 4.3 Aggregated predictions by ensemble learning

91 The aggregated predictions (mean and standard deviation) are the weighted
 92 summary of all trained models, where the weighting is the square of each model's R^2 .

$$93 \quad m_f(s, t) = \sum_i h_i(d_b, f_r) w_i \quad (S5)$$

$$94 \quad \sigma_f(s, t) = \sqrt{\frac{\sum_i w_i (h_i(d_b, f_r) - m_f(s, t))^2}{\frac{M-1}{M} \sum_i w_i}} \quad (S6)$$

$$95 \quad w_i = R_i^2 / \sum_i R_i^2 \quad (S7)$$

96 where $m_f(s, t)$ is the aggregated prediction (the weighted mean), $h_i(d_b, f_r)$ is the
 97 prediction by the i^{th} spatiotemporal model (eq. 1) trained using the bootstrap sample
 98 (d_b) and selected set of predictors (f_r); w_i is the normalized weight derived from the i^{th}

99 model's performance measure; $\sigma_f(s, t)$ is the standard deviation from the output of
100 multiple models, M is the number of nonzero weights.

101

Table S1. Correlation between average biweekly passive measurements at USC, UCLA, and UCI sites and collocated routine monitoring sites, and their linear regression coefficients used for consistent adjustment

Cover	Pollutant	Number of collocated locations	Sampling period	Correlation coefficient	Parameters		
					Slope	Intercepts	
USC	ICV1 ^a		NO ₂	Mixed dates for 2009-2013	0.94	0.89	3.57
			NO _x	Mixed dates for 2009-2013	0.93	0.83	6.23
	ICV2 ^a		NO ₂	Mixed dates for 2009-2013	0.92	0.84	3.98
			NO _x	Mixed dates for 2009-2013	0.96	0.82	5.40
UCLA	NO ₂	14	Sept. 9- Sept. 22, 2006	0.94	0.68	4.43	
			Feb. 10- Feb. 23, 2007	0.95	1.00	0.29	
	NO _x	14	Sept. 9- Sept. 22, 2006	0.98	0.80	2.38	
			Feb. 10- Feb. 23, 2007	0.97	0.81	12.05	
UCI	NO ₂	11	Jul. 10- Jul. 18, 2009	0.98	0.88	5.20	
			Jul. 24-Aug. 1, 2009	0.99	0.94	3.56	
			Nov. 13-Nov. 21, 2009	0.996	0.58	12.91	
			Dec. 4- Dec. 12, 2009	0.95	0.65	7.74	
	NO _x	11	Jul. 10- Jul. 18, 2009	0.96	0.69	8.53	
			Jul. 24-Aug. 1, 2009	0.95	0.69	5.46	
			Nov. 13-Nov. 21, 2009	0.96	1.22	-13.41	
			Dec. 4- Dec.12, 2009	0.97	0.72	14.38	

^a. ICV: The Intra-Community Variability study

Table S2. Variance explained by the predictors included in the Stage 1 mixed effects model

Covariate	Unit	Source (buffer distance)	Threshold ^a	Variances explained	
				NO ₂	NO _x
Wind speed	Meter /second	Gridded Surface Meteorological	-	3%	3%
Minimum air temperature	Celsius (°C)	Data			
		Gridded Surface Meteorological	-	5%	4%
		Data			
Spatiotemporal basis 1^b	Log ppb	Singular value decomposition by	-	19%	13%
Spatiotemporal basis 2	Log ppb	temporal basis function	-	3%	1%
CALINE4 on freeways	ppb	CALINE4 Dispersion model NO _x	>180	9%	13%
Caline4 on non-freeways	ppb	from freeway			
		CALINE4 Dispersion model NO _x	-	3%	5%
		from non-freeway			
Traffic density (300m-5km)	Vehicles/day	Distance-decayed annual traffic volume in a scaled by vehicle emission factors (in a donut radii =300 m, 5 km).	-	11%	9%
Distance to FCC1 ^c	Meter	Distance to FCC1	>15 km	7%	8%
Population Density			>21830	5%	11%
Region-level yearly mean	ppb	Annual mean concentration for the sub region determined from routine monitoring data	-	12%	13%
Spatial autocorrelation		Simulated using Thiessen polygons	-	13%	11%
Total				90%	91%

^a: threshold defined to remove the outliers for the covariate;

^b: bold font highlights the variance explained $\geq 10\%$ by the variable;

^c: Feature Class Codes for freeways and highways

Table S3. NO₂ model performance by CHS community in cross validation

ICV 1 ^a	Samples	Mean (ppb)	Correlation	RMSE ^b	NRMSE ^c	CVRMSE ^d
Alpine	156	8.71	0.89	2.51	0.1	0.29
Anaheim	128	30.07	0.69	4.02	0.13	0.13
Glendora	221	20.84	0.88	3.05	0.09	0.15
Lake Arrowhead	121	8.93	0.61	2.36	0.15	0.28
Lake Elsinore	165	10.93	0.92	1.46	0.08	0.13
Long Beach	162	20.2	0.98	2.59	0.07	0.13
Mira Loma	189	13.66	0.98	1.46	0.06	0.11
Riverside	215	16.55	0.93	2.84	0.1	0.17
San Bernardino	138	15.17	0.97	2.13	0.08	0.14
San Dimas	174	25.17	0.9	2.7	0.08	0.11
Santa Barbara	147	11.29	0.93	2.34	0.1	0.21
Santa Maria	153	8.39	0.64	1.02	0.14	0.12
Upland	232	20.14	0.95	2.73	0.08	0.14

ICV2 ^a	Samples	Mean (ppb)	Correlation	RMSE ^b	NRMSE ^c	CVRMSE ^d
Anaheim	26	19.56	0.77	3.32	0.15	0.17
Glendora	30	17.03	0.75	2.45	0.19	0.14
Long Beach	27	20.42	0.92	3.15	0.13	0.15
Mira Loma	26	18.24	0.98	1.68	0.08	0.09
Riverside	28	13.82	0.51	2.35	0.28	0.17
San Bernardino	40	11.07	0.71	2.41	0.16	0.22
San Dimas	28	20.16	0.83	2.66	0.19	0.13
Upland	28	15.79	0.75	2.56	0.18	0.16

^a: ICV, The Intra-Community Variability study; ^b: RMSE, root mean square error; ^c: RMSE, root mean square error; NRMSE, normalized RMSE; ^d:CV RMSE, coefficient of variation of the RMSE.

Table S4. NO_x model performance by CHS community in cross validation

ICV 1 ^a	Samples	Mean (ppb)	Correlation	RMSE ^b	NRMSE ^c	CVRMSE ^d
Alpine	156	18.8	0.9	4.29	0.08	0.23
Anaheim	128	67.25	0.95	6.59	0.08	0.1
Glendora	221	41.82	0.9	5.34	0.07	0.13
Lake Arrowhead	121	16.12	0.68	3.7	0.13	0.23
Lake Elsinore	165	20.34	0.91	2.5	0.08	0.12
Long Beach	162	63.01	0.98	8.78	0.07	0.14
Mira Loma	189	31.64	0.91	3.17	0.07	0.1
Riverside	215	35.86	0.91	5.99	0.08	0.17
San Bernardino	138	39.42	0.91	4.74	0.07	0.12
San Dimas	174	52.64	0.94	5.73	0.06	0.11
Santa Barbara	147	26.98	0.94	6.65	0.09	0.25
Santa Maria	153	14.21	0.67	1.43	0.13	0.1
Upland	232	41.37	0.93	5.83	0.06	0.14

ICV2 ^a	Samples	Mean (ppb)	Correlation	RMSE ^b	NRMSE ^c	CVRMSE ^d
Anaheim	26	33.4	0.76	8.86	0.17	0.27
Glendora	30	24.47	0.85	3.05	0.14	0.12
Long Beach	27	48.92	0.97	8.09	0.08	0.17
Mira Loma	26	33.9	0.93	6.98	0.12	0.21
Riverside	28	19.29	0.83	2.78	0.14	0.14
San Bernardino	40	17.58	0.83	4.21	0.11	0.24
San Dimas	28	29.51	0.83	4.43	0.18	0.15
Upland	28	22.59	0.89	3.05	0.12	0.14

^a: ICV, The Intra-Community Variability study; ^b: RMSE, root mean square error; ^c: RMSE, root mean square error; NRMSE, normalized RMSE; ^d: CV RMSE, coefficient of variation of the RMSE.

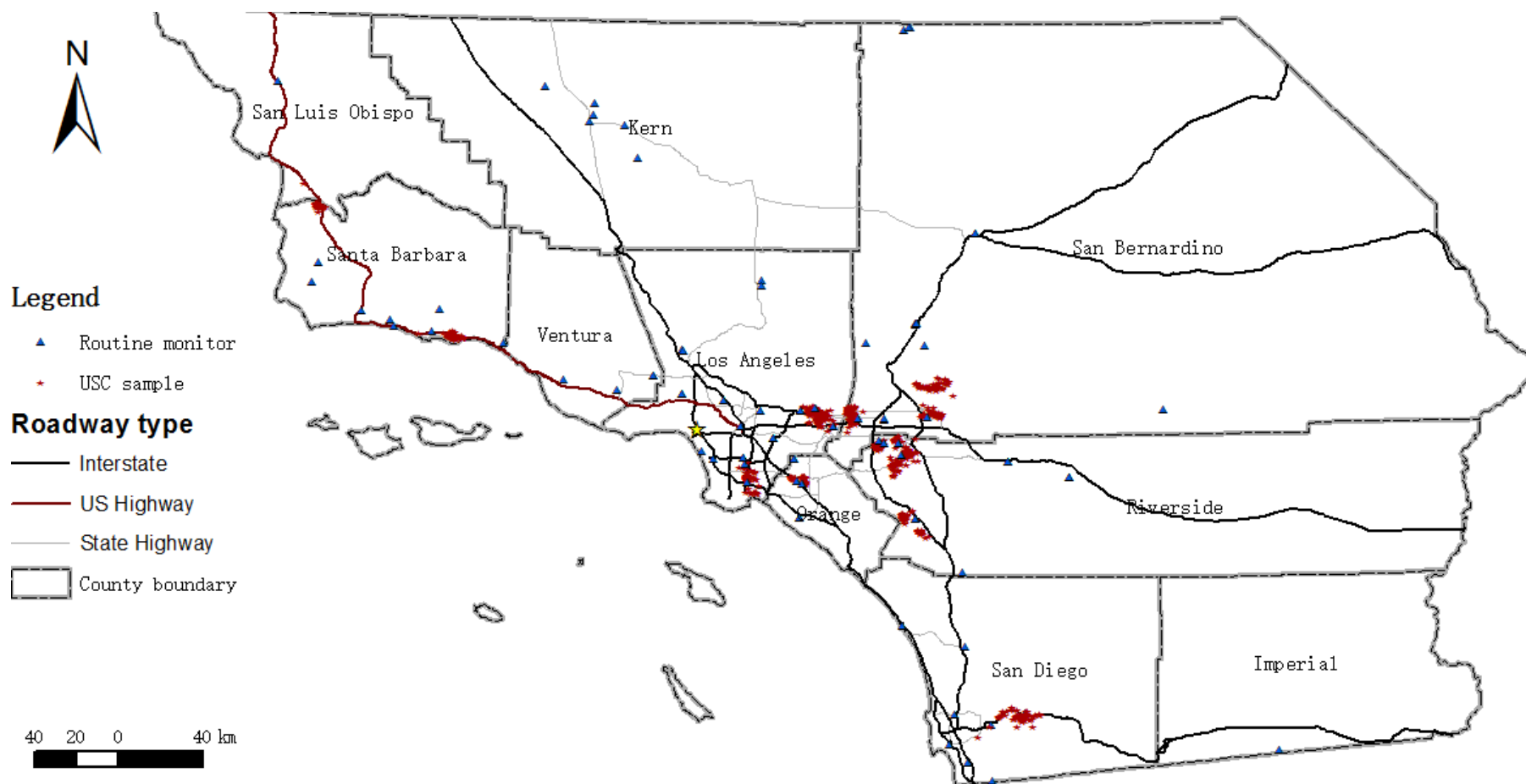
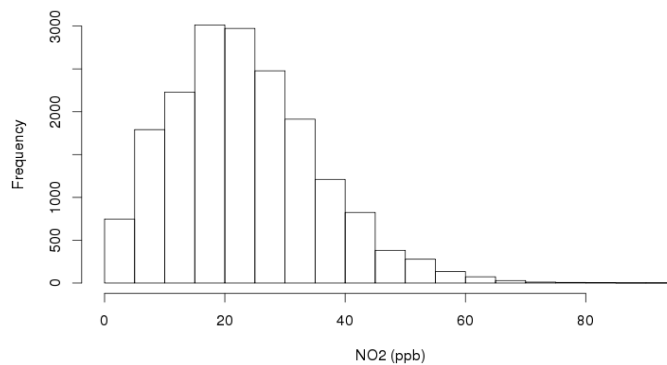
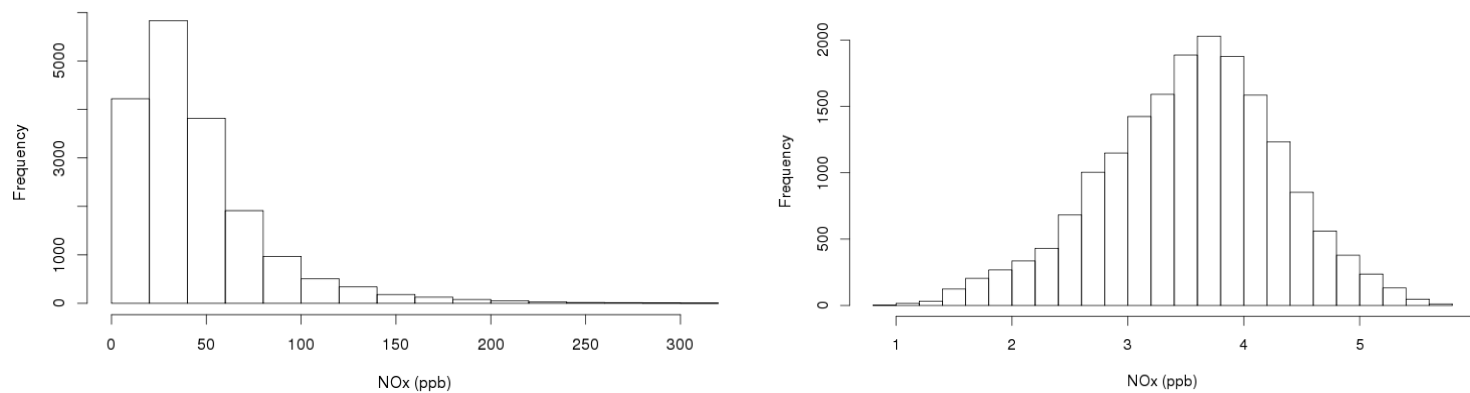


Figure S1. Study region with routine monitoring locations and USC sampling locations (the UCLA and UCI data are not shown due to IRB restrictions)

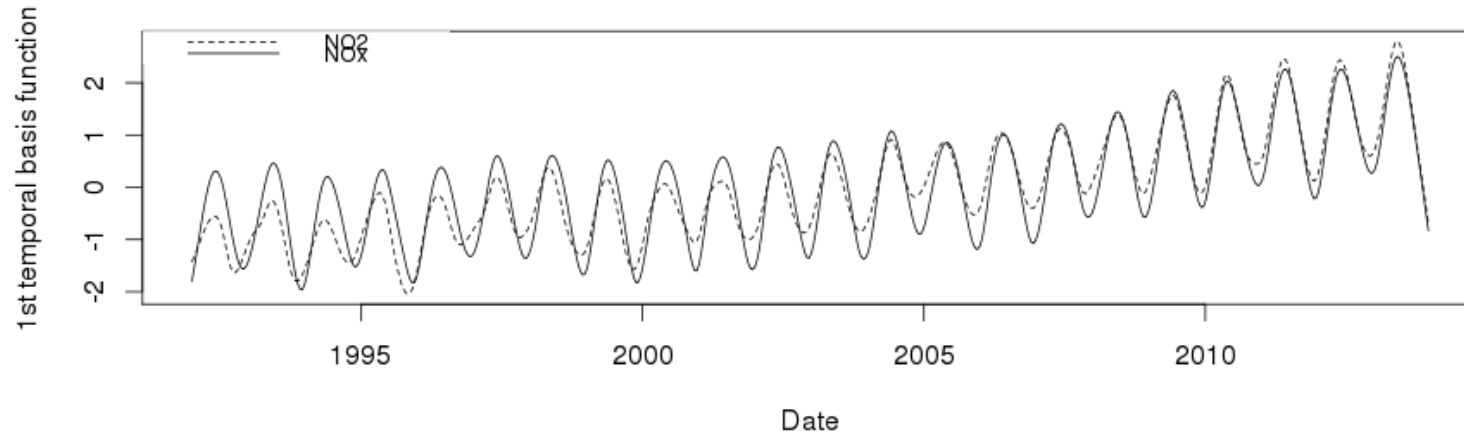


a. Histogram of NO_2 (with a small right skewness of 0.6)

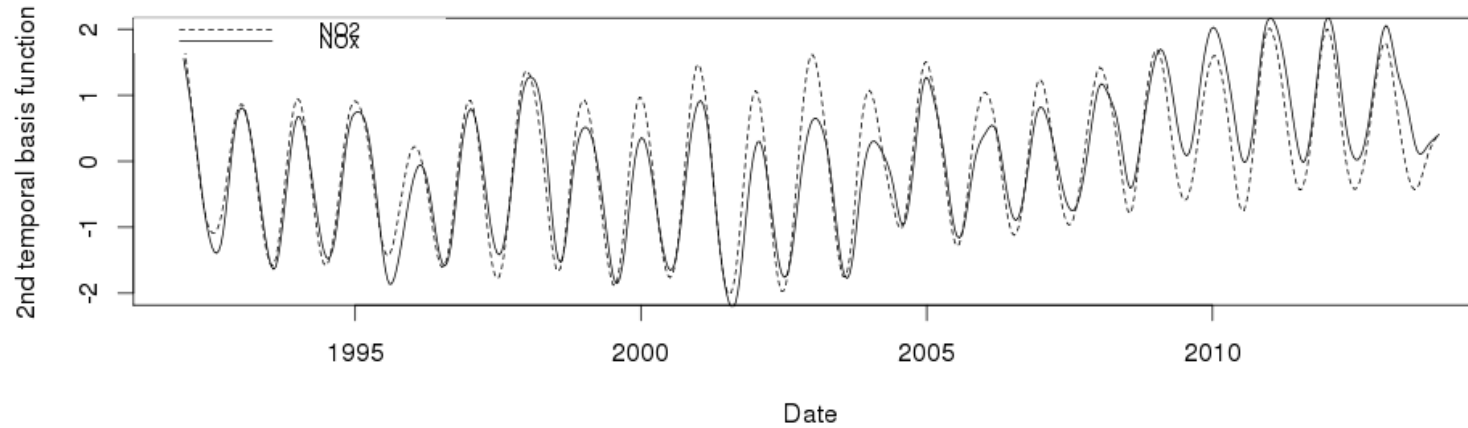


b. Histogram of NO_x (left: the original NO_x measurements with a big right skewness of 2.0; right: log-transformation with a small skewness of -0.25)

Figure S2. Histograms for NO_2 and NO_x to determine log-transformation (no log transformation for NO_2 ; log transformation for NO_x)

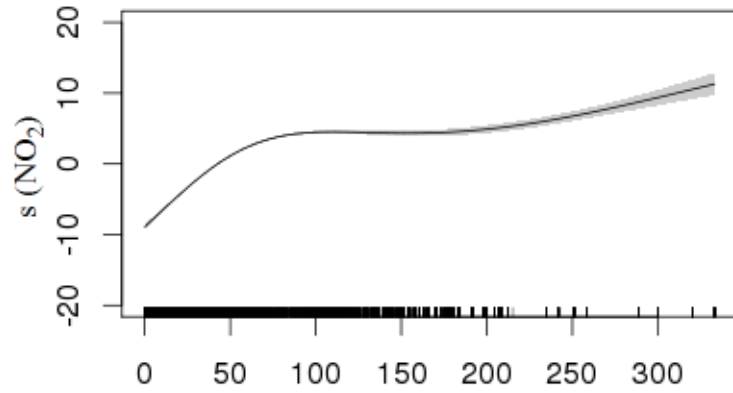


a. the first temporal basis function

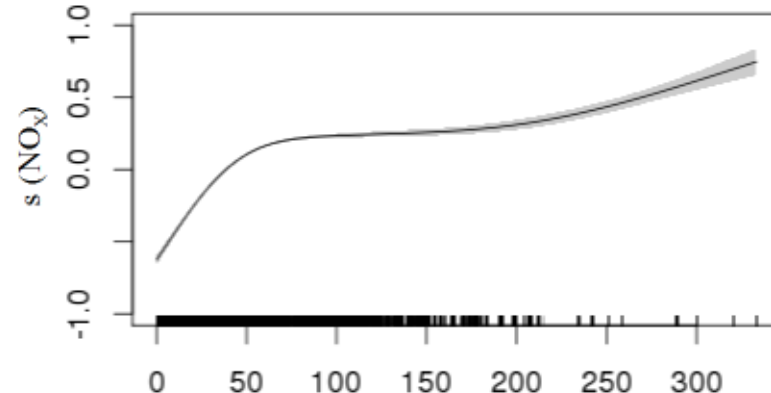


b. the second temporal basis function

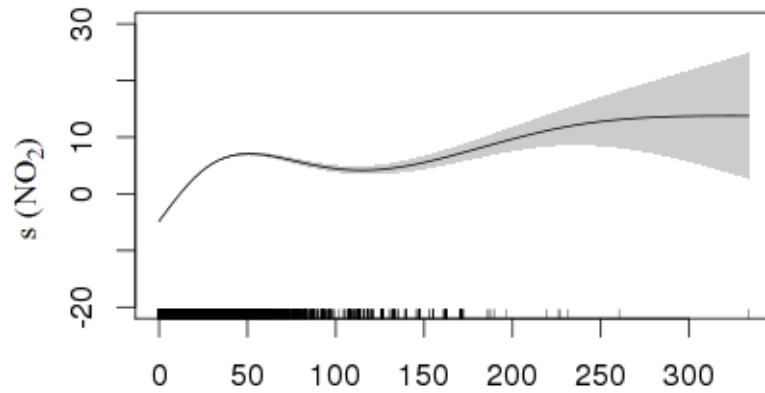
Figure S3. The first and second temporal basis functions to reflect the seasonal variability for the study region



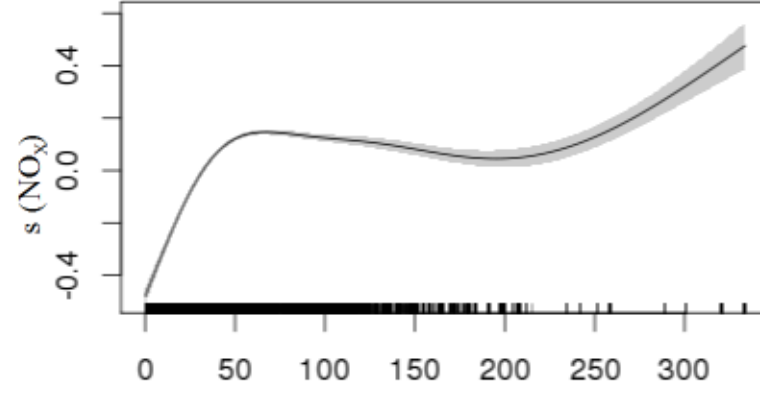
a. Traffic density (buffer distance: 0.3-5 km) for NO_2



b. Traffic density (buffer distance: 0.3-5 km) for NO_x



c. CALINE4 output (ppb) on freeways for NO_2



d. CALINE4 output (ppb) on freeways for NO_x

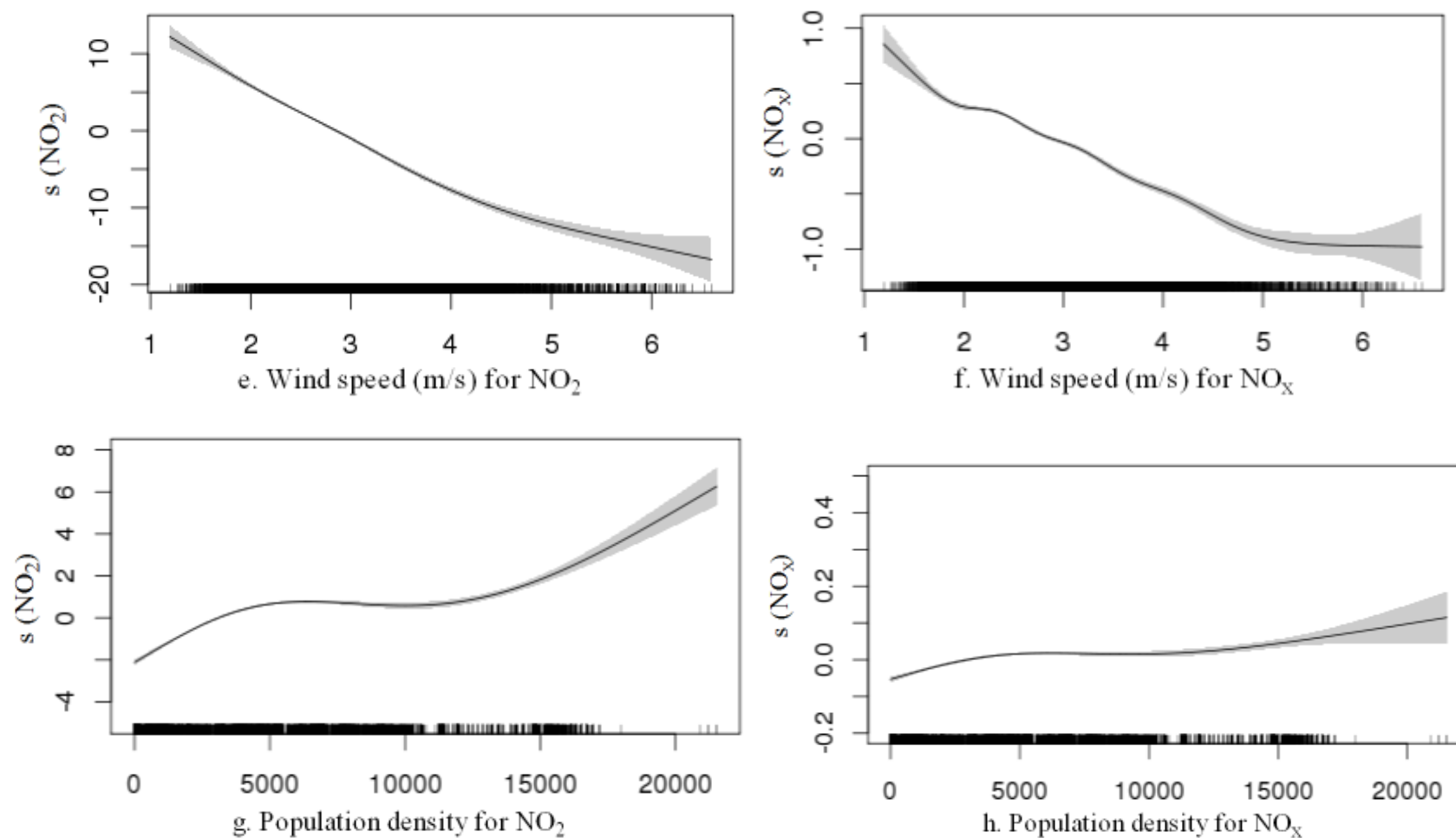


Figure S4. Non-linear association between predictive variables and concentrations by mixed models

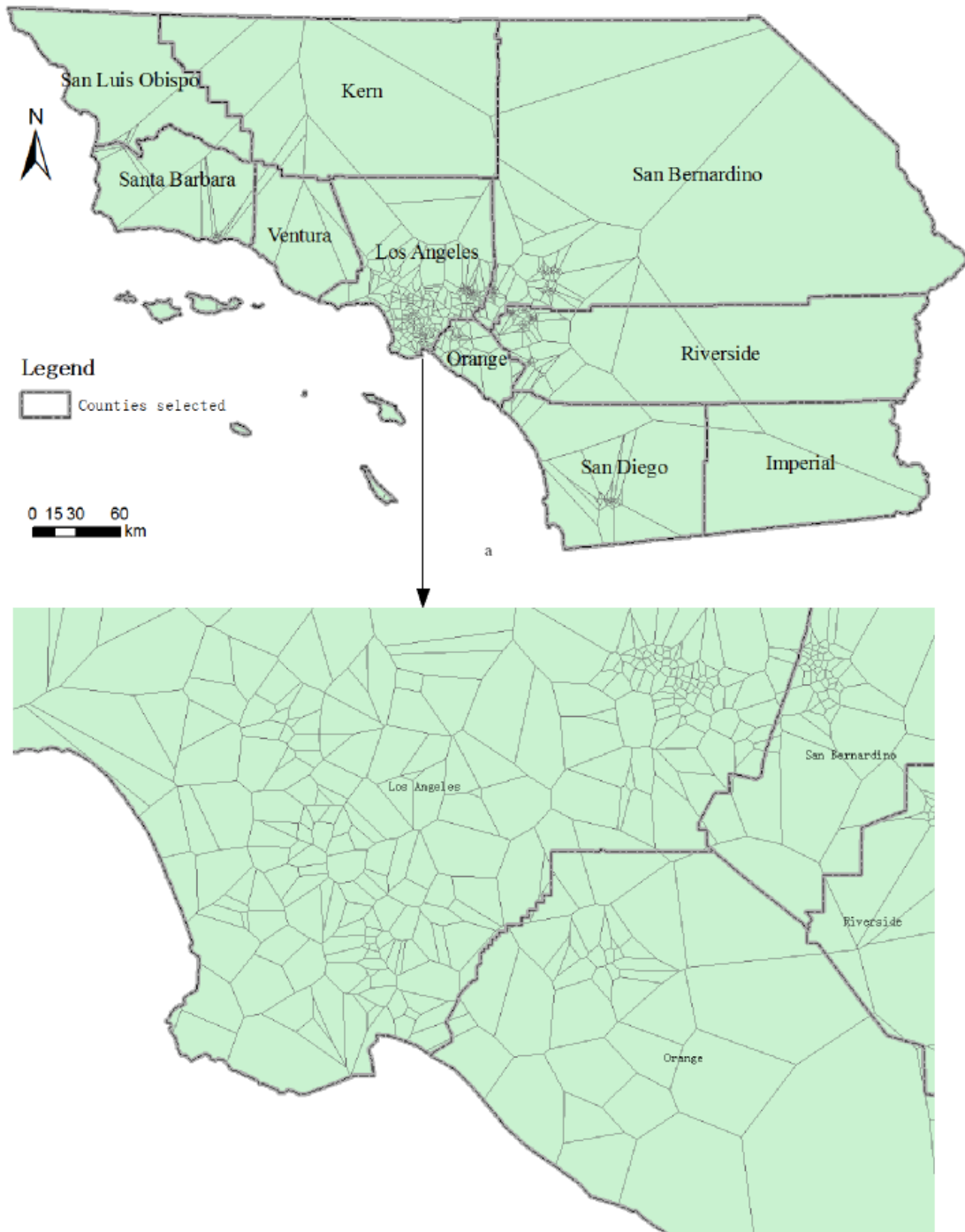
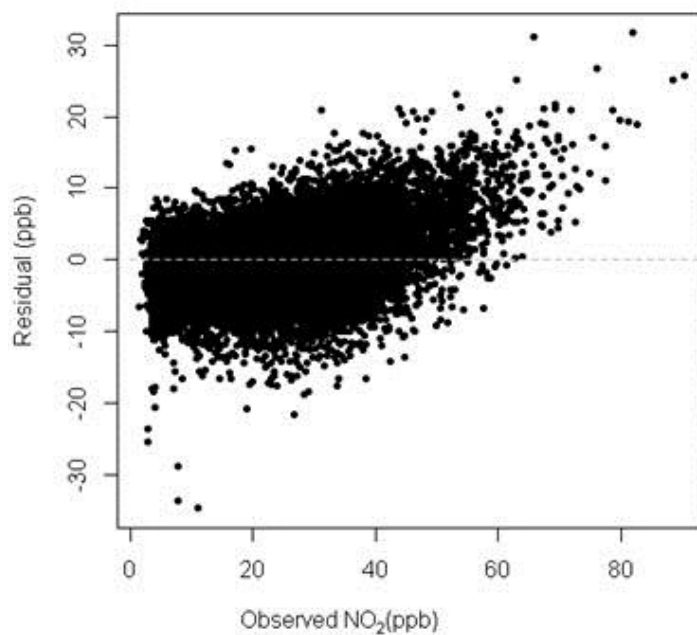
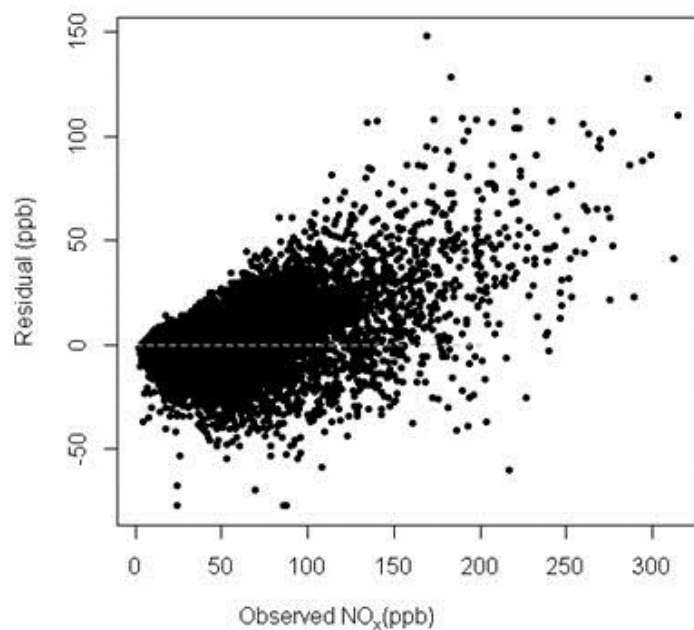


Figure S5. Spatial topology for spatial effect modeling by Thiessen polygons (aggregate distance: 500 m)

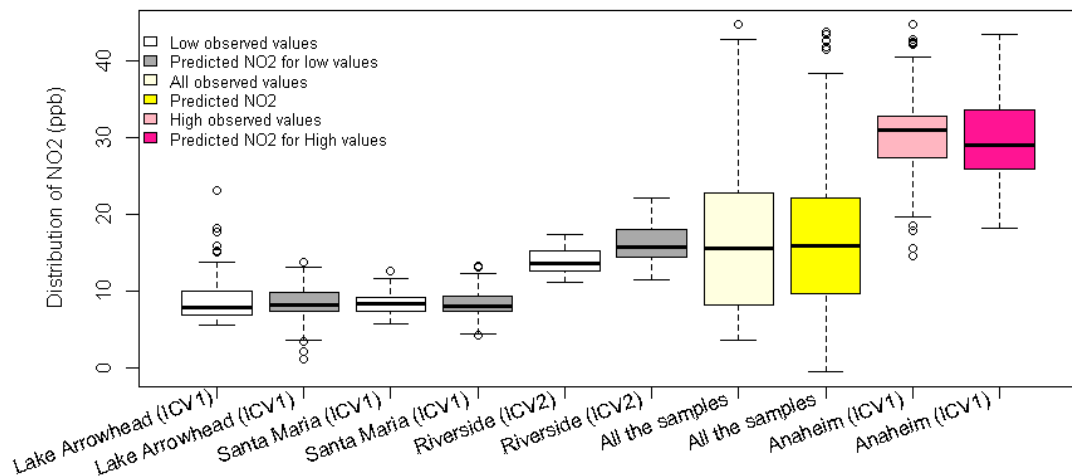


a

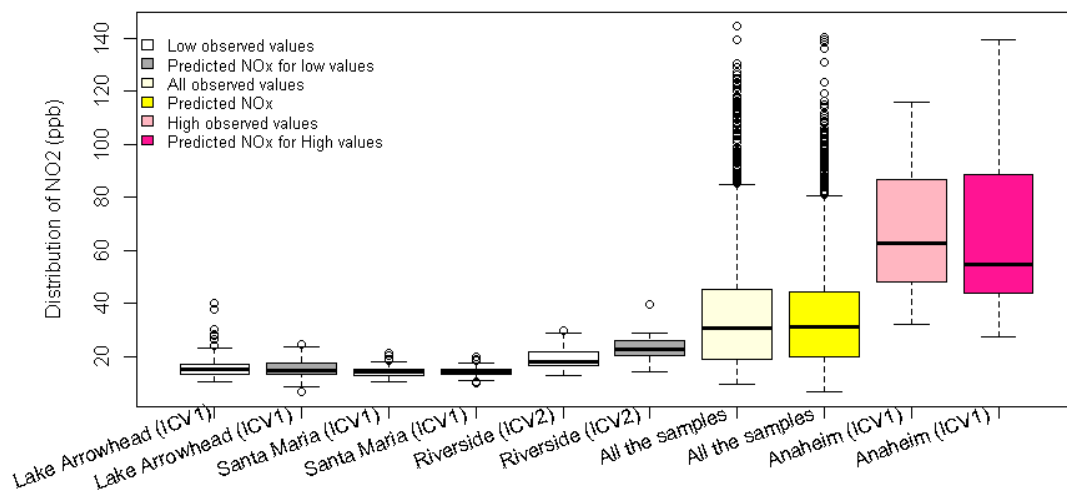


b

Figure S6. Residual plot between observed values vs. residuals (a. NO_2 ; b. NO_x) for the sample selected by bootstrap aggregating

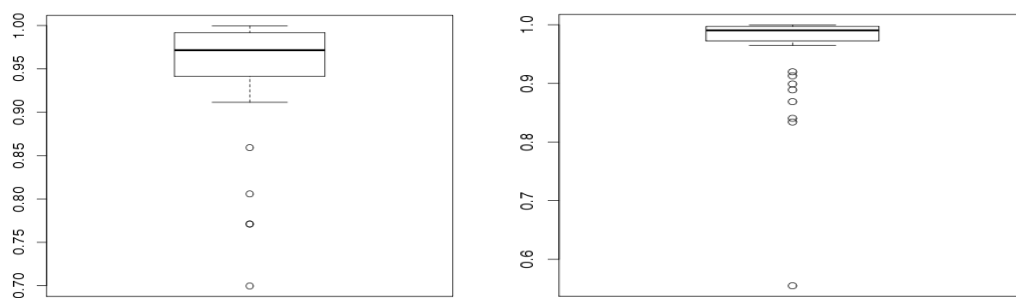


a. NO₂



b. NO_x

Figure S7. Spatial distribution of the predicted and observed NO₂ and NO_x concentration means across the CHS communities



a. NO₂

b. NO_x

Figure S8. Boxplot for correlation between constrained prediction and observed values for the time series of 51 routine monitoring stations

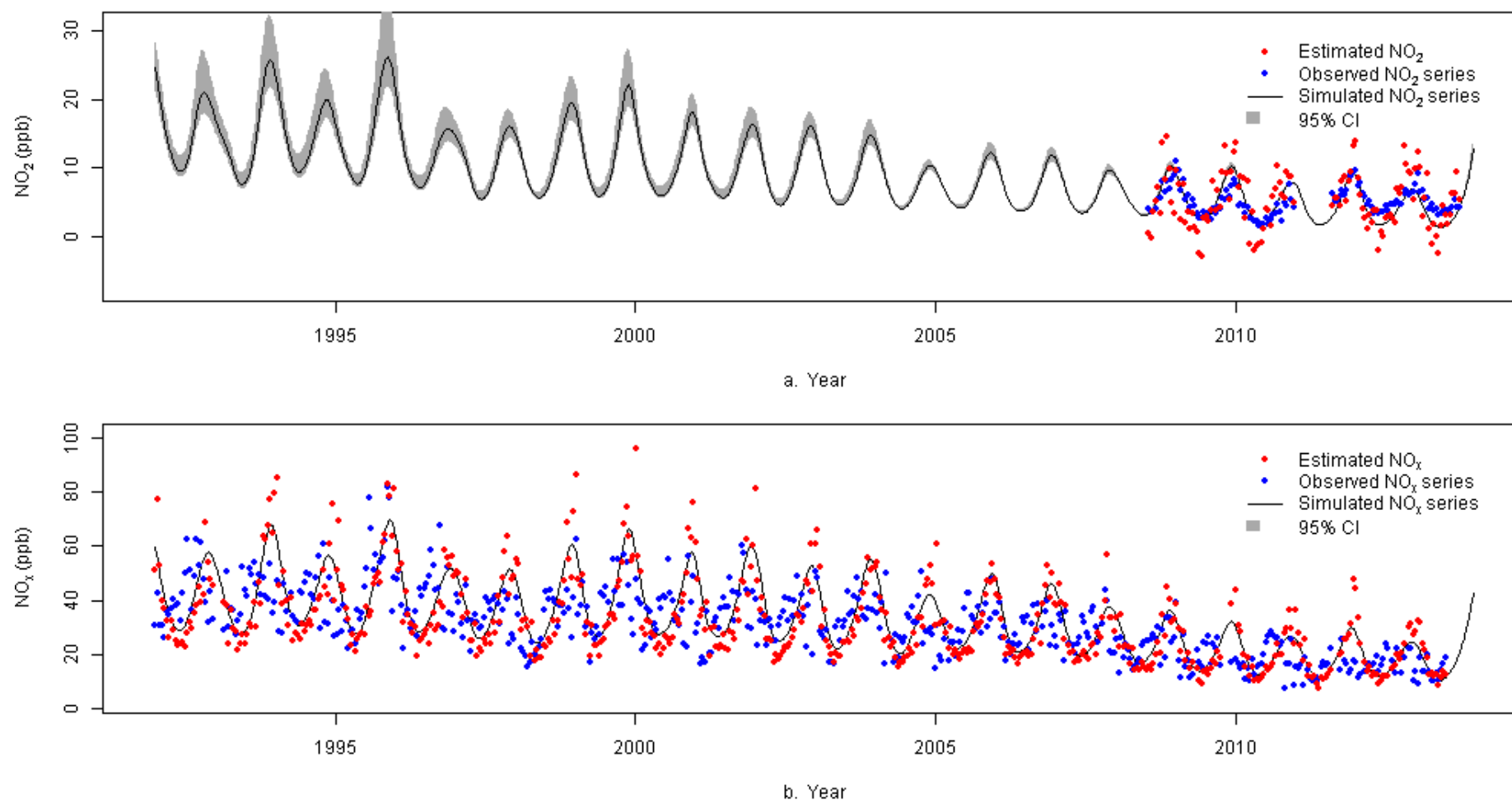
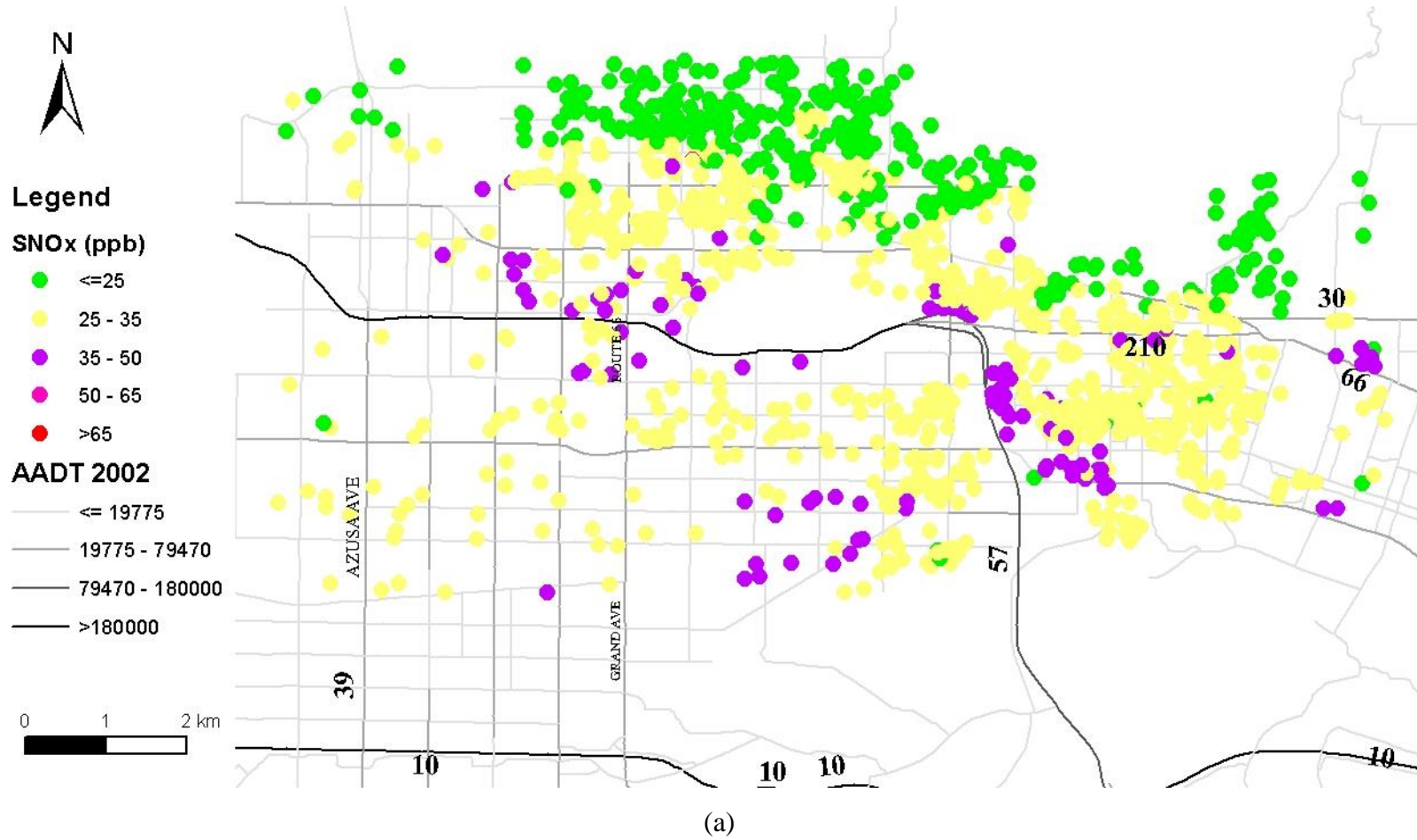


Figure S9. Time series simulated for the routine monitoring sites with the minimum correlation between constrained prediction and observed values (a.0.55 for NO_2 ; 0.70 for NO_x)



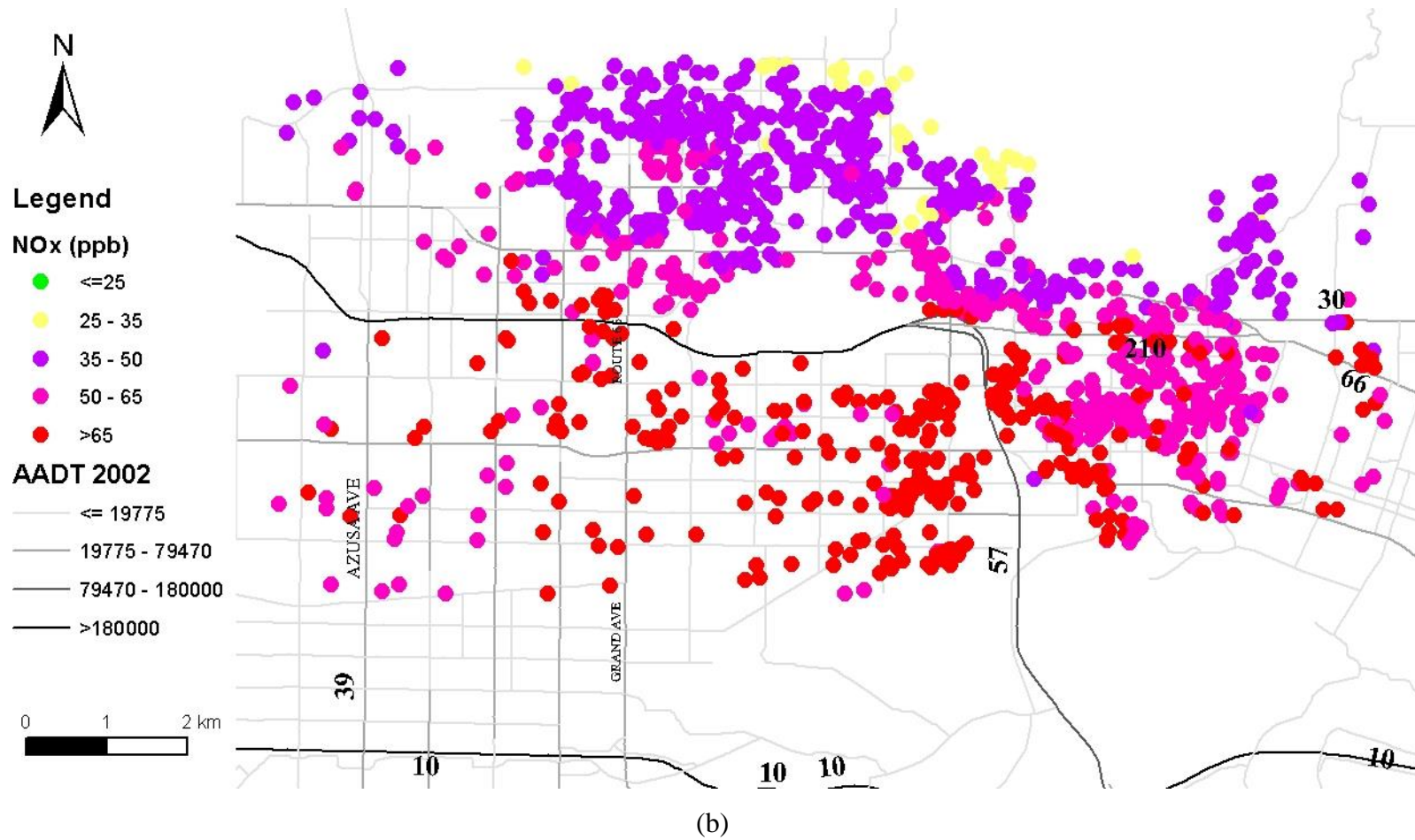


Figure S10. Summer (a) and winter (b) averages of the 2005-2006 biweekly NO_x at USC ICV1 sampling locations for San Dimas

References

1. Koutrakis, P.; Wolfson, J. M.; Bunyaviroch, A.; Froehlich, S. A. *Passive Ozone Sampler Based On A Reaction with Nitrite*; Health Effects Institute: Boston, MA, 1994; pp 19-47.
2. Franklin, M.; Vora, H.; Avol, E.; McConnell, R.; Lurmann, F.; Liu, F.; Penfold, B.; Berhane, K.; Gilliland, F.; Gauderman, W. J., Predictors of intra-community variation in air quality. *Journal of exposure science & environmental epidemiology* **2012**, *22*, (2), 135-47.
3. Su, G. J.; Jerrett, M.; Beckerman, B.; Wilhelm, M.; Ghosh, K. J.; Ritz, B., Predicting traffic-related air pollution in Los Angeles using a distance decay regression selection strategy *Environmental Research* **2009**, *109*, (2009), 657-670.
4. Li, L.; Wu, J.; Wilhelm, M.; Ritz, B., Use of generalized additive models and cokriging of spatial residuals to improve land-use regression estimates of nitrogen oxides in Southern California. *Atmospheric Environment* **2012**, *55*, 220-228.
5. Benson, P. *CALINE4: A Dispersion Model for Predicting Air Pollutant Concentrations near Roadways*; California Department of Transportation: Sacramento, CA, 1989.
6. Wu, J.; Houston, D.; Lurmann, F.; Ong, P.; Winer, A., Exposure of PM_{2.5} and EC from diesel and gasoline vehicles in communities near the Ports of Los Angeles and Long Beach, California. *Atmos Environ* **2009**, *43*, (12), 1962-1971.
7. California Air Resources Board *EMFAC2011 Technical Document*; 2011.
8. California Air Resources Board *EMFAC2014 Volume III -Technical Documentation*; 2015.
9. Besag, J.; York, J.; Mollié, A., Bayesian image restoration with two applications in spatial statistics *Annals of the Institute of Statistical Mathematics* **1991**, *43*, 1-59.