# Advances in spatial causal inference

#### **Brian Reich**

## North Carolina State University

University of Minnesota, 2020

#### Collaborators

#### Alex Larsen

#### Yawen Guan



Others: Maggie Johnson, Ana Rappold, Shu Yang and others

#### Overview

- Ambient air pollution has well-established adverse health effects
- The EPA has been regulating air pollution for many years and it has substantially reduced in the US
- The classic approach relies on a sparse network of stationary monitors
- We'll discuss methods based on two new data streams:
  - Project 1: Causal inference using numerical models to adjust for confounders (Larsen et al, in revision)
  - Project 2: Hyper-local spatiotemporal modeling using mobile monitors (Guan et al, 2020, JASA)

#### Increase in wildfires in the US

- Since the 70's, the rate of large wildfires (1000+ acres) has doubled
- In that time, the rate of very large wildfires (10,000+ acres) has increased fivefold<sup>1</sup>
- This poses an increasing health threat:
  - NPR (9/11/17) Is All That Wildfire Smoke Damaging My Lungs?
  - NYT (9/17/17): As Wildfires Burn West, Ash Rides Wind High Across the U.S.; Large Erratic Blazes are Posing a Bigger Threat to People

<sup>1</sup> climatecentral.org, 2012

#### Smoke plumes carry pollution across the continent



A B > A B
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A

#### Wildfires' contribution to air pollution

- ► Fine particular matter (PM<sub>2.5</sub>) is a criteria pollutant monitored by the EPA to protect human health
- In the US, emissions of PM<sub>2.5</sub> from most sources are steadily declining
- Forest fire smoke remains a major contributor, and may increase with climate change, population, and land-use change
- Many studies are now investigating the health effects of fire smoke

#### Causal analysis using numerical models

- ▶ How much PM<sub>2.5</sub> is causally-attributed to wildfires?
- How much health burden is causally-attributed to wildfires?
- These questions are difficult to answer because only total PM<sub>2.5</sub> (background + fire) can be measured
- Observational data, long-range transport and confounding environmental factors are also challenges
- We combine numerical models and observational data
- Causal analysis lays bare the assumptions needed to interpret the results causally

・ロト ・ 同ト ・ ヨト ・ ヨト

#### Related work on spatial causal inference

- Downscaler methods (e.g., Berrocal et al<sup>2</sup>) correct for spatially-varying bias in model output – no treatment effects
- Zigler et al have a series of papers<sup>3</sup> on causal effects of air pollution regulation – no numerical models or counterfactuals
- Detection and attribution climate studies<sup>4</sup> run models in different regimes – only observe data in one regime

<sup>4</sup>e.g., Katzfuss et al, JGR, 2017

・ 同 ト ・ ヨ ト ・ ヨ ト

<sup>&</sup>lt;sup>2</sup>e.g., Berrocal, Gelfand, and Holland, AOAS, 2010

<sup>&</sup>lt;sup>3</sup>e.g., Zigler, Dominici, Wang, Biostatistics, 2012; Zigler and Dominici, AJE, 2014; Zigler et al, HEI, 2016

## CMAQ: The Community Multiscale Air Quality Modeling System



\* plume rise, biogenic, lightning generated NO, sea salt, windblown dust, bidirectional exchange of ammonia

ж

イロト 不得 とくほと くほとう

#### Annual average CMAQ (12km $\times$ 12km)



• [1.85,3.07] • (3.07,4.4] • (4.4,6.79] • (6.79,30.4]

イロン イボン イヨン イヨン 三日

#### CMAQ run without fires



• [1.16,2.21] • (2.21,3.74] • (3.74,5.93] • (5.93,29.3]

#### Difference between the runs (as % of total)



PM<sub>2.5</sub> (%)

- (23.5,91.8]
- (12,23.5]
- (7.43,12]
- [0.84,7.43]

イロト イポト イヨト イヨト

ж

### EPA Monitoring Stations (background + fire)





- (11.5,16.8]
- (9.95,11.5]
- (7.76,9.95]
- [3.69,7.76]

イロト イポト イヨト イヨト

 $PM_{2.5}$  is measured every 3-6 days; this is the 2008-2012 average

#### Time series plot for one site in Northern CA



Brian Reich, NC State

#### Data sources and notation

- Monitor data at location s and day t:  $Y_t(s)$
- CMAQ non-fire run:  $\hat{\theta}_t(s)$
- CMAQ fire run minus non-fire run:  $\hat{\delta}_t(s)$
- Other confounders (emissions, wind, land type): X<sub>t</sub>(s)
- Binary indicator of a fire: A<sub>t</sub>(s)
- The collection a process across space is bold, e.g.,

$$A_t = \{A_t(s); s \in D\}$$

▲□▶ ▲圖▶ ▲直▶ ▲直▶ 三直 - わんで

#### Potential outcomes framework

The PM<sub>2.5</sub> at s depends on the fire status at all sites, A

This is called interference or spill-over

When envisioning counterfactual outcomes we must consider all sites simultaneously

#### Potential outcomes framework

- The "treatment" is the regime
  - R = 0: world without forest fires
  - R = 1: current world with forest fires

- $Y_t(s, 0)$  and  $Y_t(s, 1)$  are the potential PM<sub>2.5</sub> outcomes
- Model:  $Y_t(s, 0) = \theta_t(s)$  and  $Y_t(s, 1) = \theta_t(s) + \delta_t(s)$
- $\theta_t$  and  $\delta_t$  are stochastic processes

▲□▶ ▲圖▶ ▲直▶ ▲直▶ 三直 - わんで

#### Potential outcomes framework

The causal effect is

$$\Delta(\mathbf{s}) = \mathsf{E}[Y_t(\mathbf{s}, 1) - Y_t(\mathbf{s}, 0)] = \mathsf{E}[\delta_t(\mathbf{s})]$$

where the average is over the distribution of covariates X and fires A over the entire spatial domain

- Challenge: we never observe data under R = 0
- To address this we use
  - CMAQ output
  - Causal assumptions

・ 同 ト ・ ヨ ト ・ ヨ ト …

#### Assumptions

We assume there exist

- $C_t(s) \in \{0, 1\}$  where s is affected by smoke iff  $C_t(s) = 1$
- Bias-correction functions B<sub>0</sub> and B<sub>1</sub>

so that the following assumptions hold:

(A1) Consistency: 
$$Y_t(s) = Y_t[s, C_t(s)]$$

(A2) No unmeasured confounders given model output:

$$\theta_t(s) = B_0[\hat{\theta}_t(s)] + e_{1t}(s) \text{ and } \delta_t(s) = B_1[\hat{\delta}_t(s)] + e_{2t}(s),$$

where  $e_t(s) = [e_{1t}, e_{2t}]$  is independent of X, A and C

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ○ ○ ○

#### Assumptions

(A1) Consistency:  $Y_t(s) = Y_t[s, C_t(s)]$ 

(A2) No unmeasured confounders given model output:

 $\theta_t(\mathbf{s}) = B_0[\hat{\theta}_t(\mathbf{s})] + e_{1t}(\mathbf{s}) \text{ and } \delta_t(\mathbf{s}) = B_1[\hat{\delta}_t(\mathbf{s})] + e_{2t}(\mathbf{s}),$ 

where  $e_t(s) = [e_{1t}, e_{2t}]$  is independent of X, A and C

Are these assumptions reasonable?

- (A1) assumes that we have some observations we are sure are not affected by fire smoke (including spillover)...probably OK?
- (A2) assumes that the CMAQ modelers have included the important drivers of fine particulate matter....maybe OK? Have we accounted for all feedbacks?

◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

#### Causal interpretation

**Theorem 1**:Under (A1) and (A2) and assuming that  $C_t$  is not degenerate, then the parameters in the potential outcome model are identifiable via the distribution of observed data,  $Y_t \mid \hat{\theta}_t, \hat{\delta}_t, C_t$ .

- By Theorem 1, causal parameter estimation only requires inspecting the implied model for Y<sub>t</sub>(s) and confirming parameter identification.
- ► We specify parametric models for the bias correction functions B<sub>0</sub> and B<sub>1</sub> and the spatial process e<sub>t</sub>(s).
- We then argue that all parameters, including the correlation between counterfactuals, are identifiable.
- This serves as a basis for using a Bayesian approach to estimating Δ(s).

イロト 不得 とくき とくき とうき

#### Spillover/interference

- Defining the intervention as the fire regime instead of individual fires is key for two reasons:
  - 1. mimics the numerical model simulation
  - 2. limits the number of potential outcomes
- An alternative is to define potential outcomes for each potential fire-presence vector, A<sub>t</sub> = [A<sub>t</sub>(s<sub>1</sub>), ..., A<sub>t</sub>(s<sub>n</sub>)]
- This would require 2<sup>n</sup> potential outcomes, and very complicated modeling and marginalization
- Limitation: the proposed framework cannot estimate effects of individual fires or specific features of fires

#### Bayesian hierarchical model

The model is:

$$\begin{aligned} Y_t(\mathbf{s}) &= \theta_t(\mathbf{s}) + C_t(\mathbf{s})\delta_t(\mathbf{s}) + \varepsilon_t(\mathbf{s})\\ \theta_t(\mathbf{s}) &= B_0[\hat{\theta}_t(\mathbf{s})] + e_{1t}(\mathbf{s})\\ \delta_t(\mathbf{s}) &= B_1[\hat{\delta}_t(\mathbf{s})] + e_{2t}(\mathbf{s}) \end{aligned}$$

- Measurement error (iid) is  $\varepsilon$
- Background PM<sub>2.5</sub> bias adjustment:

$$B_0[\hat{\theta}_t(s)] = a_0(s) + b_0(s)\hat{\theta}_t(s)$$

Fire-contributed PM<sub>2.5</sub> bias adjustment:

$$B_1[\hat{\theta}_t(s)] = a_1(s) + b_1(s)\hat{\delta}_t(s)$$

#### Bayesian hierarchical model

The model is:

$$\begin{aligned} Y_t(\mathbf{s}) &= \theta_t(\mathbf{s}) + C_t(\mathbf{s})\delta_t(\mathbf{s}) + \varepsilon_t(\mathbf{s})\\ \theta_t(\mathbf{s}) &= B_0[\hat{\theta}_t(\mathbf{s})] + e_{1t}(\mathbf{s})\\ \delta_t(\mathbf{s}) &= B_1[\hat{\delta}_t(\mathbf{s})] + e_{2t}(\mathbf{s}) \end{aligned}$$

Fire indicator: 
$$C_t(s) = 1$$
 if  $\hat{\delta}_t(s) > \tau$ 

Spatial process e<sub>t</sub>(s) follows a bivariate Matern process

$$\operatorname{Cov}[e_{jt}(s), e_{kt}(s')] = M(||s - s'||; \rho_{jk})$$

where *M* is the Matern covariance function and  $\rho_{jk}$  are covariance parameters

#### Identification

The mean is E[Y<sub>t</sub>(s)] is

 $\alpha_0(\mathbf{s}) + \beta_0(\mathbf{s})\hat{\theta}_t(\mathbf{s}) + \alpha_1(\mathbf{s})C_t(\mathbf{s}) + \beta_1(\mathbf{s})[C_t(\mathbf{s})\hat{\delta}_t(\mathbf{s})]$ 

- ► Identification follows assuming {1, θ̂t(s), Ct(s), Ct(s), δt(s)} are not perfectly collinear
- ► The covariance Cov[Y<sub>t</sub>(s), Y<sub>t</sub>(s')] is

$$\begin{cases} M(||\mathbf{s} - \mathbf{s}'||; \rho_{11}) & \text{if } C_t(\mathbf{s}) = C_t(\mathbf{s}') = 0\\ M(||\mathbf{s} - \mathbf{s}'||; \rho_{12}) & \text{if } C_t(\mathbf{s}) \neq C_t(\mathbf{s}')\\ M(||\mathbf{s} - \mathbf{s}'||; \rho_{22}) & \text{if } C_t(\mathbf{s}) = C_t(\mathbf{s}') = 1 \end{cases}$$

By examining the correlation between pairs separately by C<sub>t</sub>(s) the correlation parameters are identifiable

#### Posterior inference

• We approximate the causal effect  $E[C_t(s)\delta_t(s)]$  by

$$\Delta(\mathbf{s}) = \sum_{t=1}^{T} C_t(\mathbf{s}) \delta_t(\mathbf{s})$$

- The Bayesian framework gives the full posterior for Δ
- Under (A1) and (A2) this has a causal interpretation
- The model is fit separately in 9 subregions
- We assume independence across days

・ロト ・ 御 ト ・ 注 ト ・ 注 ト ・

#### Data and estimates for one site in CA



Brian Reich, NC State

#### Average $\theta_t(s)$ over the 2008-2012 wildfire season



[3.55,5.41] • (5.41,6.98] • (6.98,9.82] • (9.82,15.8]
PM<sub>2.5</sub> (μg/m<sup>3</sup>)

イロン イボン イヨン イヨン 三日

#### Posterior mean slope, $b_1(s)$



• [-0.374,0.109] • (0.109,0.338] • (0.338,0.511] • (0.511,2.1]  $PM_{2.5} (\mu g/m^3)$ 

<ロ> (四) (四) (三) (三) (三)

#### Causal estimate, $\Delta(s)$ , posterior mean



#### Causal estimate, $\Delta(s)$ , posterior SD



#### Fire-Contributed PM as % of Total



#### Causal Estimate vs. CMAQ at monitors in the NW



# Estimated<sup>5</sup> number of attributable hospitalizations

		Age Group		
Region	Method	0-1	65-99	0-99
Central	Bayesian	91	170	396
	CMAQ	603	1134	2626
ENC	Bayesian	18	40	88
	CMAQ	92	205	436
South	Bayesian	111	202	456
	CMAQ	601	1092	2462
Southeast	Bayesian	196	395	901
	CMAQ	625	1260	2881
Southwest	Bayesian	100	75	248
	CMAQ	182	156	484

<sup>5</sup>Relative risk of wildfire-contributed PM<sub>2.5</sub> taken from Delfino et al (2009)

#### Estimated number of attributable hospitalizations

		Age Group		
Region	Method	0-1	65-99	0-99
Northeast	Bayesian	49	109	234
	CMAQ	178	427	895
Northwest	Bayesian	89	194	423
	CMAQ	177	392	848
West	Bayesian	427	714	1631
	CMAQ	912	1548	3507
WNC	Bayesian	47	28	130
	CMAQ	93	57	257

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

#### Cumulative health burden by county



イロン イボン イヨン イヨン

#### Summary and future projects

- We derived the assumptions needed to make causal claims using numerical model output
- We found large contributions of fires to PM, especially in the west, even after adjusting for model bias

#### Future work:

- Continuous treatments
- Causal quantile effects
- Analyze individual fires

# Changing scales!

- Air pollution epidemiology relies on a few stationary monitors per city
- The field is undergoing a paradigm shift due to fine-resolution mobile monitors
- We analyze data from sensors strapped to Google StreetView cars
- Some cities now have thousands of low-cost sensors
- Phone apps are under development



Photo: Apte (2017)

・ 同 ト ・ ヨ ト ・ ヨ ト

#### Mobile data: Oakland NO<sub>2</sub>

- Two cars were deployed from June 2015 to May 2016
- Starts on weekdays at  $\approx$  9am, drove  $\approx$  6-8 hours each day
- Measurements are taken at roughly every second.
- Large missing data (car maintenance, sensor failure, etc.)

#### Statistical Challenges:

- ► Data are large (≈900,000 observations in current dataset)
  - Computationally prohibitive to fit full spatiotemporal models
- Data are extremely sparse in space and time
  - Maximum of two observations at a time point.
  - Only a small region of Oakland is sampled on each day.
- Data are noisy and subject to outliers

# Example daily observations of log(NO<sub>2</sub>)



- Car A and B drove from 8am 2pm
- 12,389 observations, covered less than a third of Oakland

#### Objectives

Develop a statistical model incorporating landuse covariates and spatiotemporal dependence for real-time, high resolution (30m) forecasting of air pollution.

- 1. How well and how far ahead can we reasonably forecast air pollution levels?
- 2. If we were to design a new Google car study, how many cars should be deployed to improve prediction?
- 3. Is deploying sensors on cars more efficient than a fixed-location sensor network?

# Temporal aggregation



We took temporal block medians to dampen effects of extremes

• • • • • • • • • •

< ∃⇒

#### Example landuse covariates



Highway
Major
Residential

Commerical Industrial Residential NA

#### PCA of landuse variables



#### Non-spatial landuse regression

Let  $Y_t(s)$  be the log(NO<sub>2</sub>) at time *t* and location s

$$Y_t(\mathbf{s}) = \mathsf{X}_t(\mathbf{s})^T \boldsymbol{\beta} + \epsilon_t(\mathbf{s}), \quad \epsilon_t(\mathbf{s}) \stackrel{\textit{iid}}{\sim} \mathcal{N}(\mathbf{0}, \tau^2)$$

where  $X_t(s)$  contains

- The first seven PCs
- Four trig functions for hourly diurnal cycle
- Interactions between the PCs and trig functions
- $R^2 \approx 0.16$  and residuals are correlated

(日本) (日本) (日本)

#### Results from landuse regression

Observed vs. Predicted, Oct. 29, 2015 - Dec. 18, 2015



< D > < P > < E</p>

#### Results from landuse regression

Observed vs. Predicted, Oct. 29, 2015 - Dec. 18, 2015



≣⇒

< < >> < </>

#### Spatiotemporal landuse regression model

We add a spatiotemporal process to capture dependence

$$\mathsf{Y}_t(\mathsf{s}) = \mathsf{X}_t(\mathsf{s})eta + rac{\eta_t(\mathsf{s})}{\eta_t} + \epsilon_t(\mathsf{s}), \quad \epsilon_t(\mathsf{s}) \stackrel{\textit{iid}}{\sim} \mathsf{N}(\mathsf{0}, au^2)$$

The covariance is

$$Cov\left[\eta_t(\mathbf{s}), \eta_{t'}(\mathbf{s}')\right] = \sigma^2 \exp\left\{-\sqrt{||\mathbf{s} - \mathbf{s}'||^2/\rho + |t - t'|^2/\phi}\right\}$$

 $\triangleright$   $\rho$  and  $\phi$  determine the spatial and temporal dependence.

(4個) (4回) (4回)

•

#### Computation

- MLE for the full dataset is impossible
- We first estimate  $\beta$  using OLS and compute residuals  $e_i$
- We use the Veccia approximation to estimate the covariance parameters
- The joint distribution is approximated by the product of conditional distributions

$$f(\boldsymbol{e}_1,...,\boldsymbol{e}_n) \approx \prod_{i=1}^n f\left(\boldsymbol{e}_i | \boldsymbol{e}_j \in \mathcal{N}_i\right)$$

where  $N_i$  is the set of "neighbors" for observation *i* 

This is fast if the neighboring sets are small

ъ

#### Computation

The neighboring sets are observations in the recent past

 $\mathcal{N}_i = \{ \text{obs between } I \text{ and } I + 60 \text{ minutes prior to obs } i \}$ 

Taking *I* = 0 gives the best approximation to the likelihood

 Gramacy and Apley (2015) show its better to include some distant neighbors

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

#### Prediction correlation using 1 sec block medians

		Prediction lag			
Model	1	5 mins	15 mins	60 mins	Car AB
Х	-	0.18	0.18	0.18	0.08
S	-	0.27	0.27	0.27	0.09
ST	0	0.45	0.25	0.18	0.09
ST	5	0.58	0.36	0.28	0.10
ST	15	0.57	0.36	0.31	0.10
ST	60	0.55	0.38	0.28	0.09

ъ

・ロト ・ 同ト ・ ヨト ・ ヨト

#### Prediction correlation using 1 min block medians

		Prediction lag			
Model	1	5 mins	15 mins	60 mins	Car AB
Х	-	0.28	0.28	0.28	0.19
S	-	0.34	0.34	0.34	0.21
ST	0	0.59	0.44	0.29	0.26
ST	5	0.64	0.56	0.46	0.26
ST	15	0.64	0.56	0.45	0.26
ST	60	0.63	0.55	0.45	0.26

3

・ロト ・ 同ト ・ ヨト ・ ヨト

#### Forecast for Dec 2015 - Feb 2016



#### 15 minutes ahead forecasts of NO<sub>2</sub>



- Forecast for 15:00 using the data from 13:45 to 14:45 on May 5, 2016
- As expected, standard errors are lowest where data has been obtained most recently from the two cars.

#### Network design

**Objective:** To improve NO<sub>2</sub> prediction how many Google cars should be deployed? How many fixed-location sensors would provide the same quality of prediction?



Deploy different number of mobile and fixed-location sensors.



< ロ > < 同 > < 三

# Prediction performance comparison.

#### Summary and future projects

- Our work shows that short-term forecasts of air pollution at a high spatial resolution are possible
- Future work:
  - Fuse mobile and stationary sensors
  - Model extremes
  - Multicity analysis
  - Design efficient routes
- Works supported by NIH and NSF

#### THANKS!