

Unmeasured spatial confounding

Georgia Papadogeorgou

University of Florida

Learning from Interventions, Simons Institute, February 17 2022

Joint work with Cory Zigler (UT Austin),
Christine Choirat (Swiss Data Science Center),
and Patrick Schnell (Ohio State)

Spatial confounding: What does it mean?

- Spatial statistics and causal inference have different views
- Common thread: Confounding **hides** what we want to learn
- Difference: The goal, **what it is** that we want to learn

Spatial confounding: What does it mean?

- Spatial statistics and causal inference have different views
- Common thread: Confounding **hides** what we want to learn
- Difference: The goal, **what it is** that we want to learn

• Notation	\rightsquigarrow	Treatment or exposure	Z_i
		Potential outcomes	$Y_i(z)$
		Outcome	$Y_i = Y_i(Z_i)$
		Measured covariates	C_i

The goal: Learn the relationship $Y \mid Z, \mathbf{C}$

- Spatial Model 1: $Y \sim Z + \mathbf{C}$
 \rightsquigarrow Residuals are spatially correlated!

The goal: Learn the relationship $Y \mid Z, \mathbf{C}$

- Spatial Model 1: $Y \sim Z + \mathbf{C}$
 - ↪ Residuals are spatially correlated!
- Spatial Model 2: $Y \sim Z + \mathbf{C} + U$
 - ↪ U spatial random effect
 - ↪ Reason for U : Inference, capture residual spatial dependence
 - ↪ U is given a correlation structure (Exponential, Matérn, etc)

The goal: Learn the relationship $Y \mid Z, \mathbf{C}$

- Spatial Model 1: $Y \sim Z + \mathbf{C}$
 - ~ Residuals are spatially correlated!
- Spatial Model 2: $Y \sim Z + \mathbf{C} + U$
 - ~ U spatial random effect
 - ~ Reason for U : Inference, capture residual spatial dependence
 - ~ U is given a correlation structure (Exponential, Matérn, etc)
- U is collinear with the exposure Z
 - ~ “Steals” from the exposure
 - ~ **Confounding by the spatial U**

Hodges and Reich (2010); Paciorek (2010); Hughes and Haran (2013); Hanks et al. (2015);
Vicente et al. (2020); Azevedo et al. (2020); Reich et al. (2020)

The goal: Learn the relationship $Y \mid Z, \mathbf{L}$

where \mathbf{L} such that $Z \perp\!\!\!\perp Y(z) \mid \mathbf{L}$

Here: $\mathbf{L} = (\mathbf{C}, U)$, for U spatial

The goal: Learn the relationship $Y \mid Z, \mathbf{L}$

where \mathbf{L} such that $Z \perp\!\!\!\perp Y(z) \mid \mathbf{L}$

Here: $\mathbf{L} = (\mathbf{C}, U)$, for U spatial

- Spatial Model 1: $Y \sim Z + \mathbf{C}$
 - ~> does not include all confounders
 - ~> cannot be used to learn causal effects

The goal: Learn the relationship $Y \mid Z, \mathbf{L}$

where \mathbf{L} such that $Z \perp\!\!\!\perp Y(z) \mid \mathbf{L}$

Here: $\mathbf{L} = (\mathbf{C}, U)$, for U spatial

- Spatial Model 1: $Y \sim Z + \mathbf{C}$
 - ↪ does not include all confounders
 - ↪ cannot be used to learn causal effects
- Spatial Model 2: $Y \sim Z + \mathbf{C} + U$
 - ↪ U spatial random effect
 - ↪ U cannot “learn” the unmeasured variable (Paciorek, 2010; Schnell and Papadogeorgou, 2020)
 - ↪ cannot be used to learn causal effects

Spatial confounding: What does it mean?

Spatial statistics

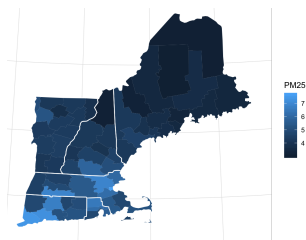
- Want to learn the relationship between outcome and exposure given the measured variables
- Spatial adjustment is for inference
- Collinearity of exposure and random effect “blurs” the results

Causal inference

- Want to learn the relationship between outcome and exposure conditioning on **all** confounders
- Some confounders are spatial and unmeasured
- Spatial models cannot adjust for unmeasured spatial confounders
- Can we use unmeasured confounders' *structure* to adjust for them?

Question 1: Spatial causal inference in air pollution research

- Regulations enforce stricter rules on emissions to reduce air pollution
 - ↪ Power plants follow various strategies to comply to these regulations
- Install SCR/SNCR systems for reducing NO_x emissions
 - ↪ NO_x reacts with VOCs and carbon monoxide in the presence of sunlight to create ozone
- The scientific questions are **causal**
 - ↪ SCR/SNCR systems VS alternatives on ambient air pollution
- The data are **spatial**
 - ↪ Exposure, outcome, measured and unmeasured covariates are spatially structured
 - ↪ VOCs, sunlight spatial & unmeasured



Question 2: Supermarket access and cardiovascular health

- Access to supermarkets \rightsquigarrow Healthy habits \rightsquigarrow Cardiovascular health
- Potential confounders: Income, demographics, regional culture, personal vehicles, diet, local-level support for individuals with disabilities
- Data

⎧	Exposure	Continuous $\in (0, 100)$
	Level	Measured at counties
	Confounders	Unmeasured, hard to define

Two approaches to unmeasured spatial confounding

- 1 “Causal” approach
Incorporating spatial information in propensity score methods
with Christine Choirat, Cory Zigler
- 2 “Spatial” approach
Bias correction in outcome regression
with Patrick Schnell

Distance adjusted propensity score matching for binary treatments

Distance Adjusted Propensity Score Matching

- Estimate the ATT = $E[Y(1) - Y(0)|Z = 1]$
- Propensity score model using measured variables C :

$$P(Z_i = 1|C_i) = \text{expit}(C_i^T \beta)$$

- For a treated unit i and a control unit j define

$$DAPS_{ij} = w|PS_i - PS_j| + (1 - w) * \text{Dist}_{ij}, \quad w \in [0, 1]$$

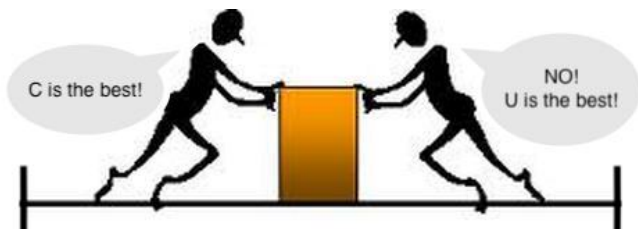
where PS propensity score estimates, and Dist spatial proximity.

- Small value $DAPS_{ij}$ means:
 - Similar propensity scores
 - Points in close geographical distance (similar values of U !)

Choosing w

w : relative importance of the observed and unobserved confounders

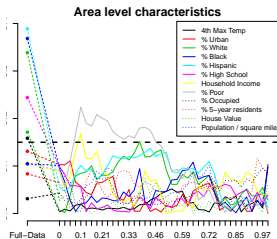
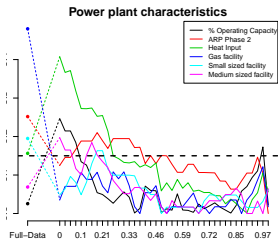
- High values of w - priority to observed covariates
- Low values of w - priority to spatial proximity



Navigate the tradeoff between:

- 1 Making matches as similar as possible with respect to C
- 2 Small distance of matched pairs to capture similarity in U

Balance and distance

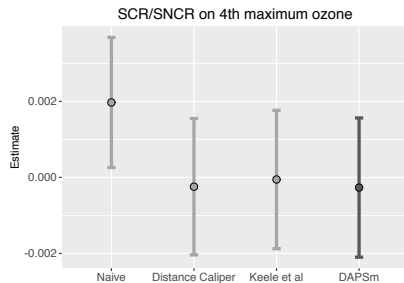
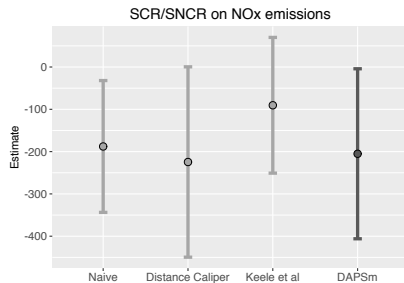


Naive pairs



DAPSm pairs





Keele et al. (2015)

Outcome regression with bias correction to mitigate bias from unmeasured spatial variables

- We assume the following *true* model for the potential outcomes:

$$Y_i(z) = \eta(z, \mathbf{C}) + g(\mathbf{W}^u) + \varepsilon_i(z)$$

- \mathbf{W}^u are unmeasured variables
- Additive model, \mathbf{W}^u do not interact with Z and \mathbf{C}

- We assume the following *true* model for the potential outcomes:

$$Y_i(z) = \eta(z, \mathbf{C}) + U + \varepsilon_i(z)$$

- \mathbf{W}^u are unmeasured variables
- Additive model, \mathbf{W}^u do not interact with Z and \mathbf{C}
- Denote $U = g(\mathbf{W}^u)$

- We assume the following *true* model for the potential outcomes:

$$Y_i(z) = \beta_0 + \beta_1 z + U + \varepsilon_i(z)$$

- \mathbf{W}^u are unmeasured variables
- Additive model, \mathbf{W}^u do not interact with Z and \mathbf{C}
- Denote $U = g(\mathbf{W}^u)$
- For ease of presentation, assume \mathbf{C} empty, $\eta(z) = \beta_0 + \beta_1 z$
- Focus on $\beta_1 = \mathbb{E}[Y(z+1) - Y(z)]$

If we could fit model $Y \sim Z + U$, we could estimate β_1 without bias.

If we could fit model $Y \sim Z + U$, we could estimate β_1 without bias.

- $Y \sim Z \rightarrow \widehat{\beta}$

- $Y \sim Z + \text{Spatial RE} \rightarrow \widetilde{\beta}$

If we could fit model $Y \sim Z + U$, we could estimate β_1 without bias.

- $Y \sim Z \rightarrow \widehat{\beta}$

- Bias = $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{E}(U|Z)$

- $Y \sim Z + \text{Spatial RE} \rightarrow \widetilde{\beta}$

- Bias = $\{\mathbf{X}^\top (\text{Var}[\mathbf{Y}|Z])^{-1} \mathbf{X}\}^{-1} \mathbf{X}^\top (\text{Var}[\mathbf{Y}|Z])^{-1} \mathbf{E}[U|Z]$

where $\mathbf{X} = (\mathbf{1}, Z)$

If we could fit model $Y \sim Z + U$, we could estimate β_1 without bias.

- $Y \sim Z \rightarrow \widehat{\beta}$

- Bias = $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{E}(U|Z)$

- $Y \sim Z + \text{Spatial RE} \rightarrow \widetilde{\beta}$

- Bias = $\{\mathbf{X}^\top (\text{Var}[\mathbf{Y}|Z])^{-1} \mathbf{X}\}^{-1} \mathbf{X}^\top (\text{Var}[\mathbf{Y}|Z])^{-1} \mathbf{E}[U|Z]$

where $\mathbf{X} = (\mathbf{1}, Z)$

- Identify the bias term, and subtract it

$$\bar{\beta} = \{\mathbf{X}^\top (\text{Var}[\mathbf{Y}|Z])^{-1} \mathbf{X}\}^{-1} \mathbf{X}^\top (\text{Var}[\mathbf{Y}|Z])^{-1} \{\mathbf{Y} - \mathbf{E}[U|Z]\}$$

- Find a way to identify $E[U|Z]$!

- ① (\mathbf{U}, \mathbf{Z}) is mean 0 normal

$$\begin{pmatrix} \mathbf{U} \\ \mathbf{Z} \end{pmatrix} \sim \mathcal{N} \left[\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{G} & \mathbf{Q} \\ \mathbf{Q}^\top & \mathbf{H} \end{pmatrix}^{-1} \right]$$

- ② **Cross-Markov property:** $p(Z_i | \mathbf{Z}_{-i}, \mathbf{U}) = p(Z_i | \mathbf{Z}_{-i}, U_i)$
 $\leadsto Q$ is diagonal

- ③ **Constant conditional correlation:** $\text{Cor}(U_i, Z_i | \mathbf{U}_{-i}, \mathbf{Z}_{-i}) = \rho$
 $\leadsto q_{ii} = -\rho \sqrt{g_{ii} h_{ii}}$

Calculating the affine estimator

- Integrating $\mathbf{U}|\mathbf{Z}$ out

$$\begin{aligned}\mathbf{Z} &\sim \mathcal{N}[\mathbf{0}, (\mathbf{H} - \mathbf{Q}^\top \mathbf{G}^{-1} \mathbf{Q})^{-1}] \\ \mathbf{Y}|\mathbf{Z} &\sim \mathcal{N}[\mathbf{X}\boldsymbol{\beta} - \underbrace{\mathbf{G}^{-1} \mathbf{Q} \mathbf{Z}}_{\mathbb{E}[\mathbf{U}|\mathbf{Z}]}, \mathbf{G}^{-1} + \mathbf{R}^{-1}],\end{aligned}$$

where $\mathbf{R}^{-1} = \text{Cov}(\boldsymbol{\varepsilon})$

- The restricted likelihood is

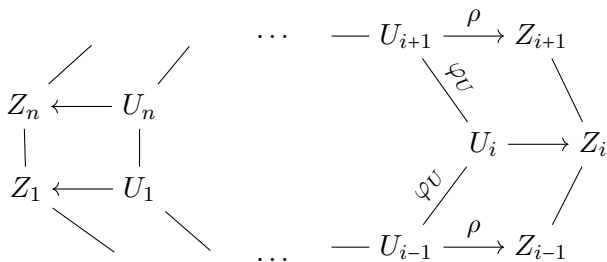
$$RL \propto \exp \left[-\frac{1}{2} \left\{ (\mathbf{Y} - \mathbf{B}\mathbf{Z})^\top \mathbf{C}_2 (\mathbf{Y} - \mathbf{B}\mathbf{Z}) + \mathbf{Z}^\top \mathbf{A}^{-1} \mathbf{Z} \right\} \right]$$

where $\mathbf{A} = (\mathbf{H} - \mathbf{Q}^\top \mathbf{G}^{-1} \mathbf{Q})^{-1}$, and $\mathbf{B} = -\mathbf{G}^{-1} \mathbf{Q}$

- We can calculate $\bar{\boldsymbol{\beta}}$ using the RL maximizers / Bayesian

Learning the spatial parameters

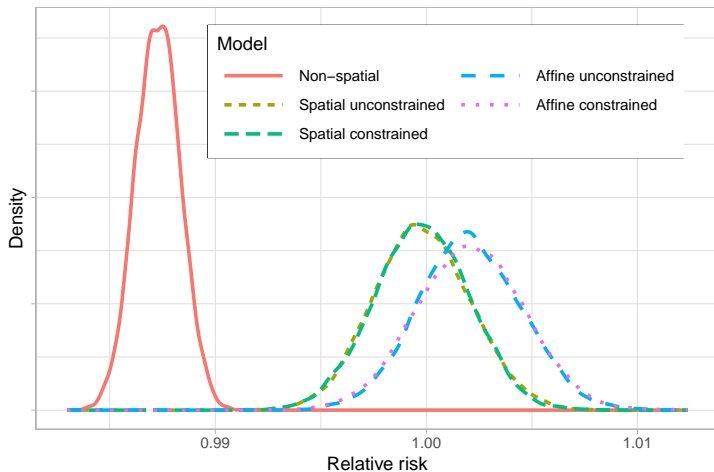
- (1) the unmeasured U is (a) spatial and (b) a confounder
- (2) the dependencies can be depicted on a ring
→ everything is identifiable based on (Z, Y)



Conjecture: Similar results hold for graphs with “enough pairs” of locations at varying lags

If U not spatial: Effect is not identifiable

Effect of poor supermarket access on CVD deaths



Conclusions

- Spatial statistics VS causality: different goals, different tools
- In causal inference with structured data it is possible to use the structure information to alleviate unmeasured confounding bias
- Identifiability is hinged on structure (spatial correlation)
- Interplay between the scale of variation in the unmeasured variable and the causal positivity assumption
→ scale restriction in estimation

}	Paper 1	Data, R package, PDF gpapadogeorgou.netlify.app/publication/dapsm/
	Paper 2	Data, Code, PDF gpapadogeorgou.netlify.app/publication/spatial_confounding2/

- D. R. M. Azevedo, M. O. Prates, and D. Bandyopadhyay. Alleviating spatial confounding in spatial frailty models. *arXiv:2008.06911*, 2020.
- E. M. Hanks, E. M. Schliep, M. B. Hooten, and J. A. Hoeting. Restricted spatial regression in practice: geostatistical models, confounding, and robustness under model misspecification. *Environmetrics*, 26(4):243–254, 2015.
- J. S. Hodges and B. J. Reich. Adding Spatially-Correlated Errors Can Mess Up the Fixed Effect You Love. *The American Statistician*, 64(4):325–334, 2010.
- J. Hughes and M. Haran. Dimension reduction and alleviation of confounding for spatial generalized linear mixed models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(1):139–159, 2013.
- L. Keele, R. Titiunik, and J. Zubizarreta. Enhancing a Geographic Regression Discontinuity Design Through Matching to Estimate the Effect of Ballot Initiatives on Voter Turnout. *Journal of Royal Statistical Society A*, 178:223–239, 2015.
- C. J. Paciorek. The importance of scale for spatial-confounding bias and precision of spatial regression estimators. *Statistical Science*, 25(1):107, 2010.
- G. Papadogeorgou, C. Choirat, and C. M. Zigler. Adjusting for unmeasured spatial confounding with distance adjusted propensity score matching. *Biostatistics*, 20(2):256–272, 2019.
- B. J. Reich, S. Yang, Y. Guan, A. B. Giffin, M. J. Miller, and A. G. Rappold. A review of spatial causal inference methods for environmental and epidemiological applications. *arXiv:2007.02714*, 2020.
- P. Schnell and G. Papadogeorgou. Mitigating unobserved spatial confounding when estimating the effect of supermarket access on cardiovascular disease deaths. *Annals of Applied Statistics*, 14(4):2069–2095, 2020.
- G. Vicente, T. Goicoa, P. Fernandez-Rasines, and M. Ugarte. Crime against women in india: unveiling spatial patterns and temporal trends of dowry deaths in the districts of uttar pradesh. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 183(2): 655–679, 2020.