



Modern Spatial Path Analytic Tools to Investigate the Geography of Medical Debt across a US State

Modern Modeling Methods – 2024, Storrs CT,
June 25-26, 2024 Session 1B: Modeling Spatial Data

Slides at <https://tinyurl.com/mmmdebtct>



Emil Coman, PStat, SEMNET ‘moderator’; Samuel Bruder
comanus@gmail.com

General plan

1. Causal logic with spatial data
2. Spatial non-independence intuition
 - Modeling solutions
3. Example with $N=8$ regions

'Visual' causal reasoning

Figure S2· Directed acyclic graph showing selected factors involved in the lifetime risk of major adverse cardiovascular events (MACE) after childhood cancer survivorship

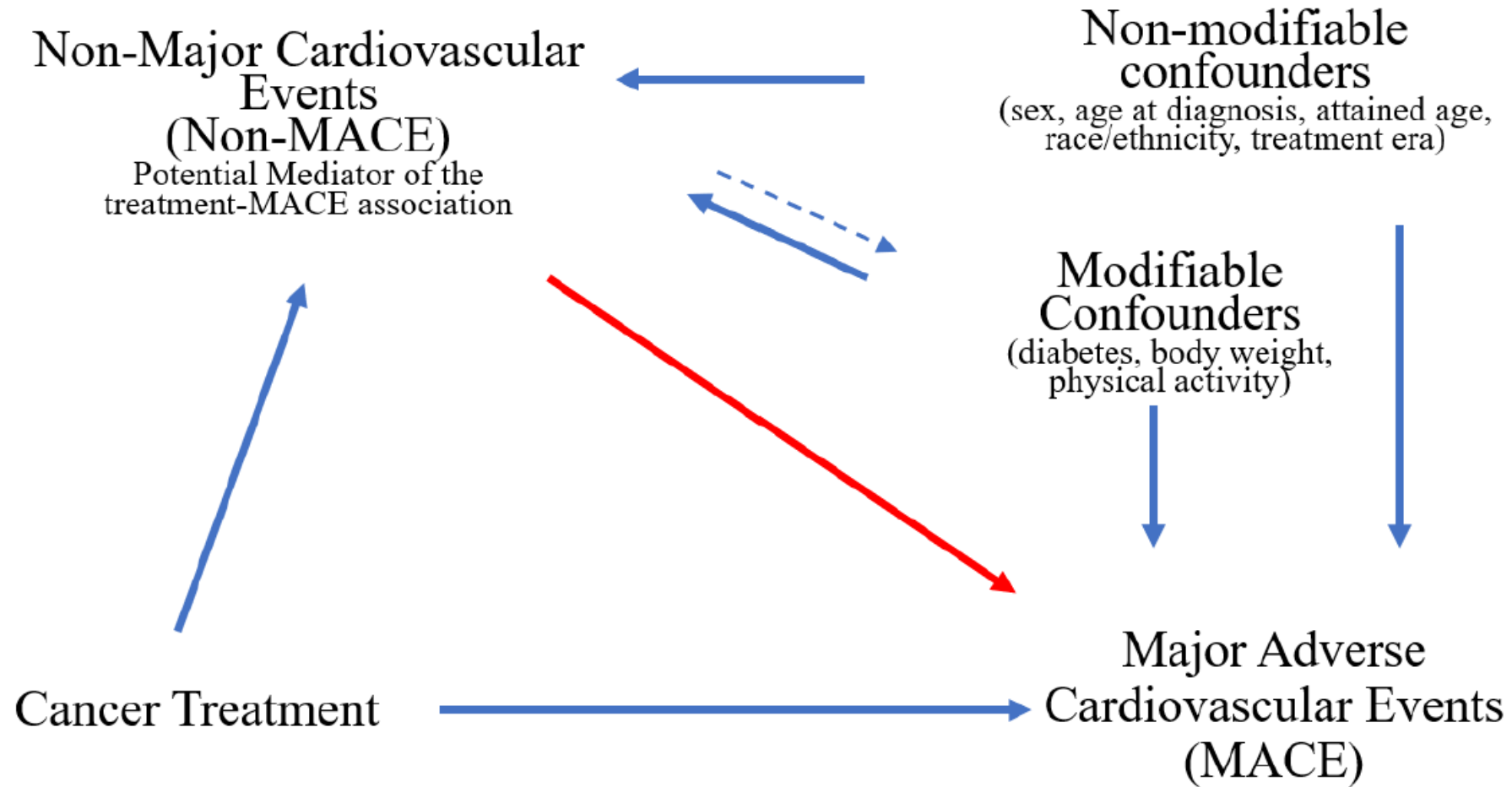
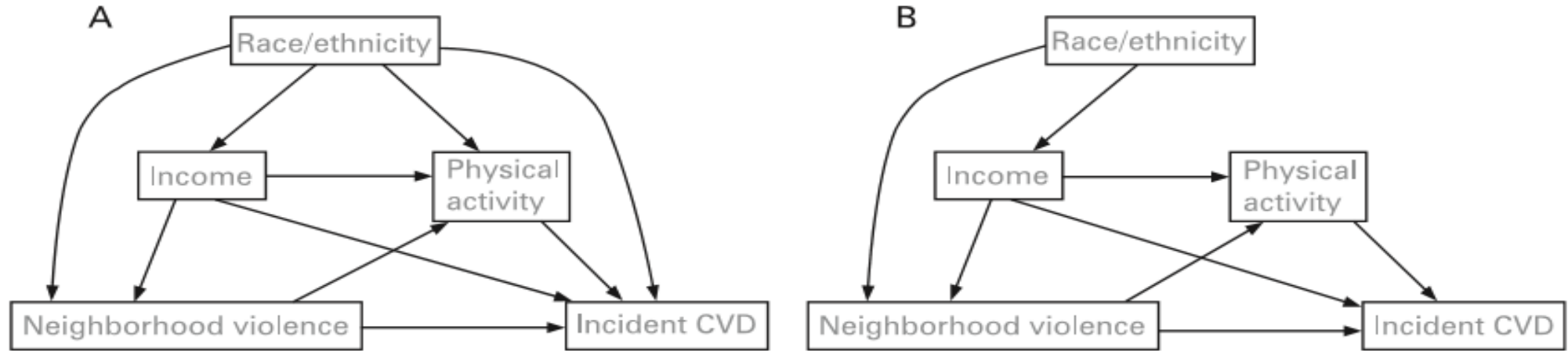
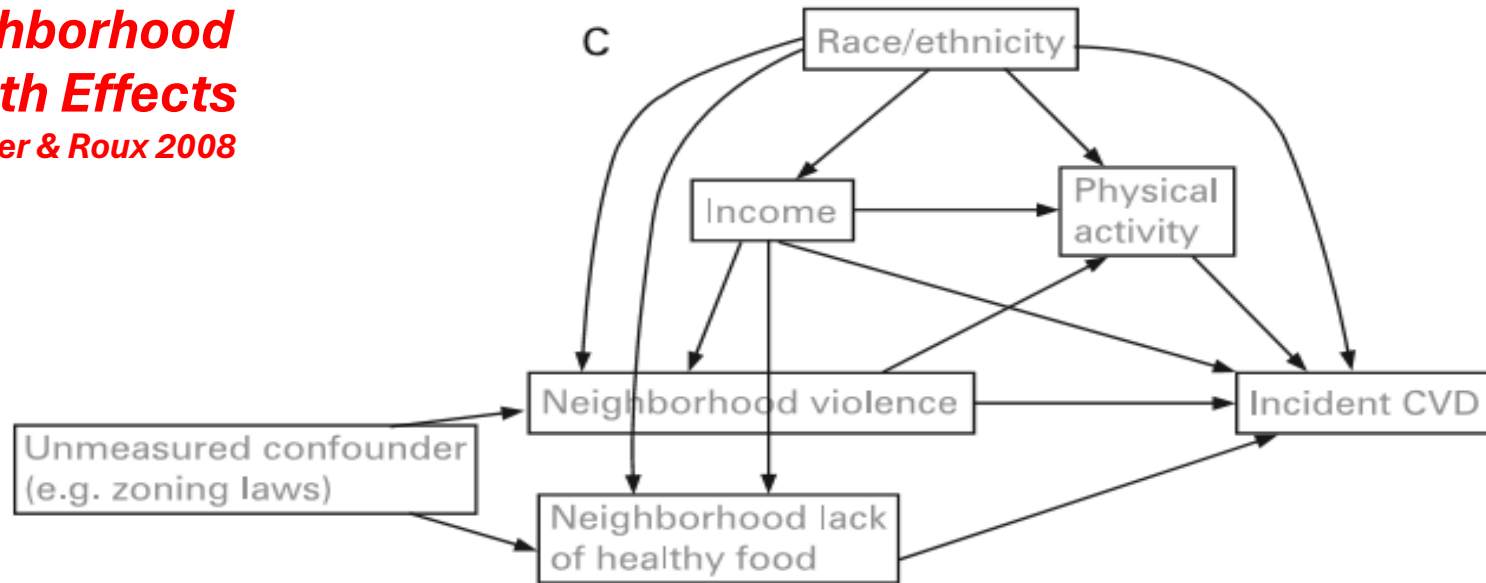


Figure 1 Using directed acyclic graphs to identify variables that need to be controlled for in estimating neighbourhood health effects.



**Neighborhood
Health Effects**
Fleischer & Roux 2008



Causal reasoning

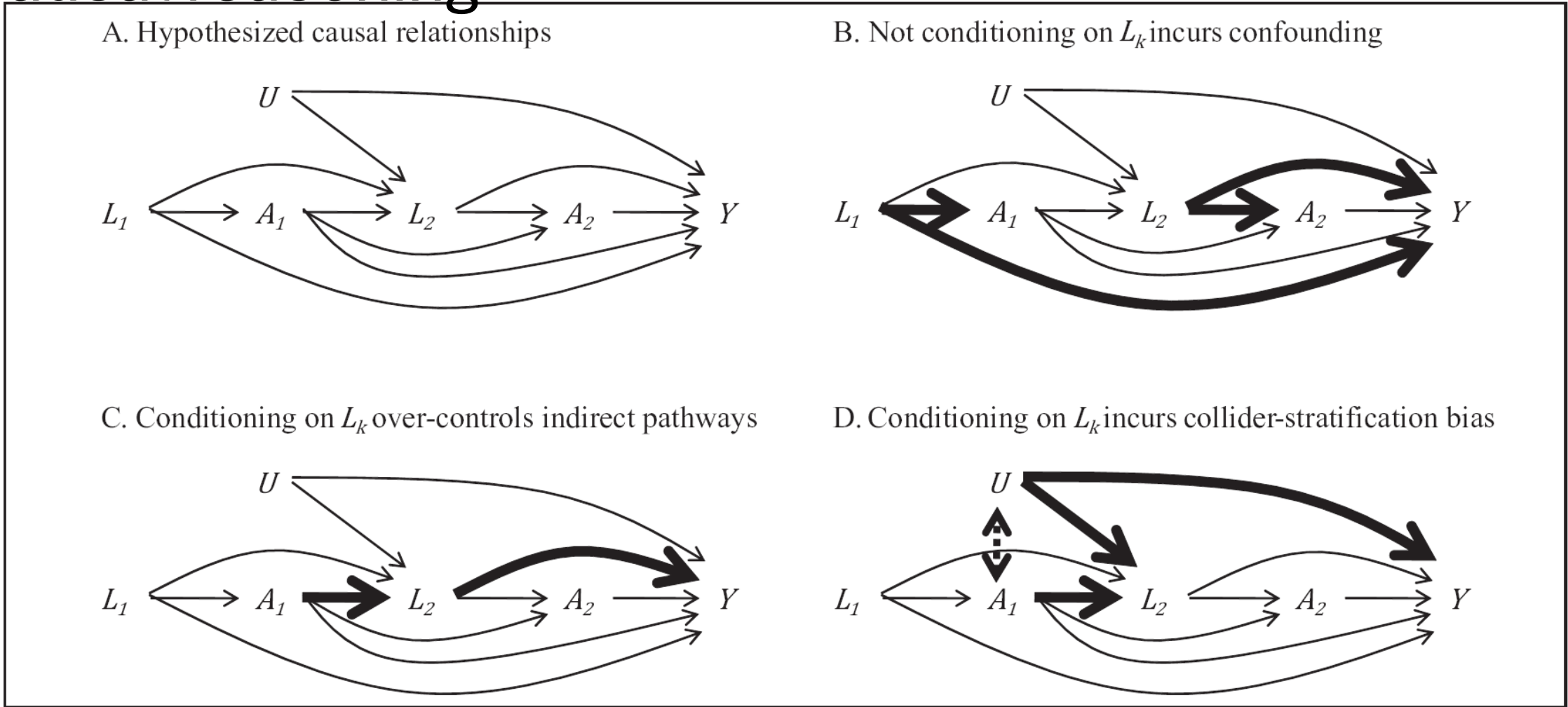
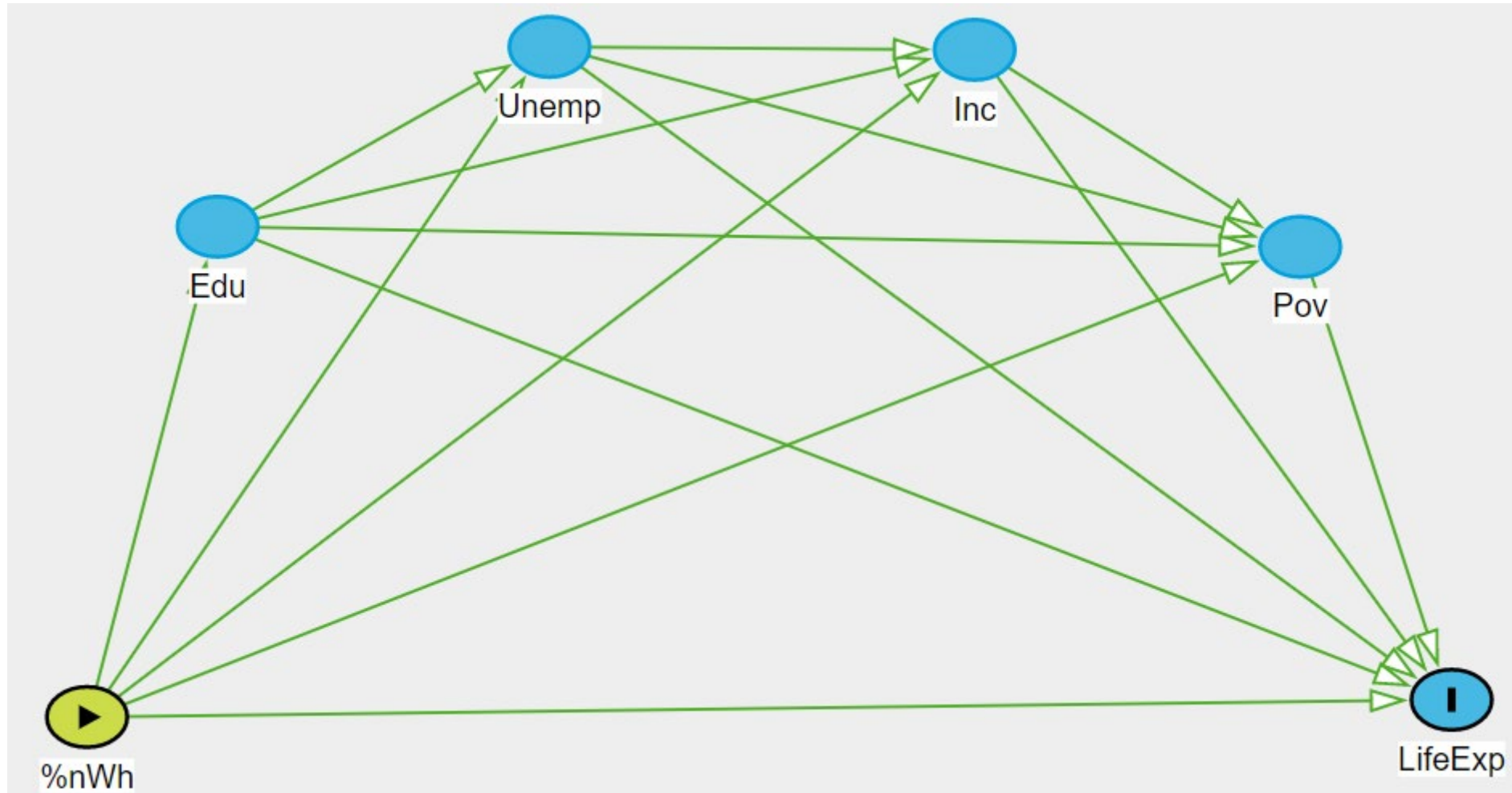


Figure 1. Causal Graphs for Exposure to Disadvantaged Neighborhoods with Two Waves of Follow-up

Note: A_k = neighborhood context, L_k = observed time-varying confounders, U = unobserved factors, Y = outcome.

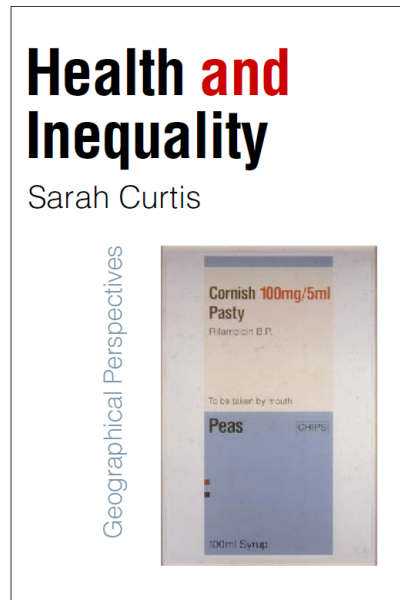
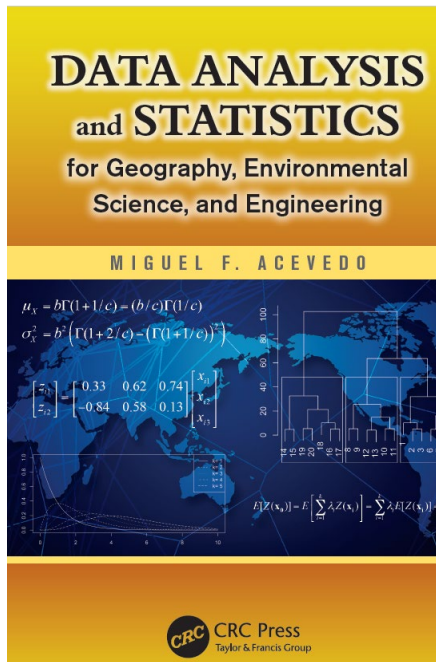
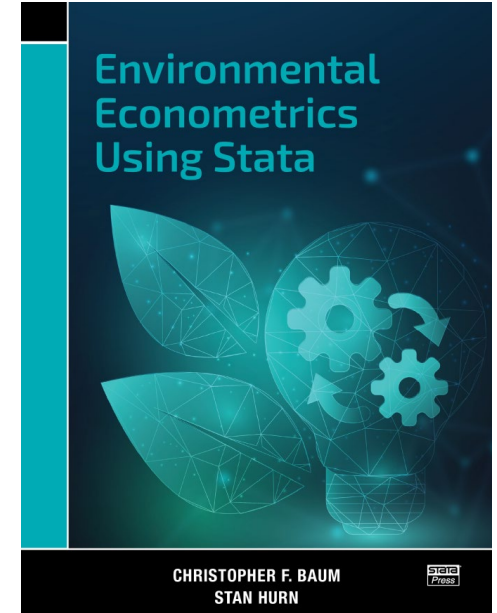
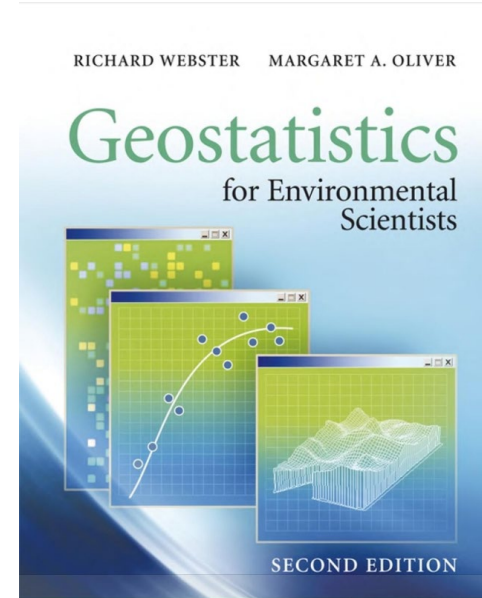
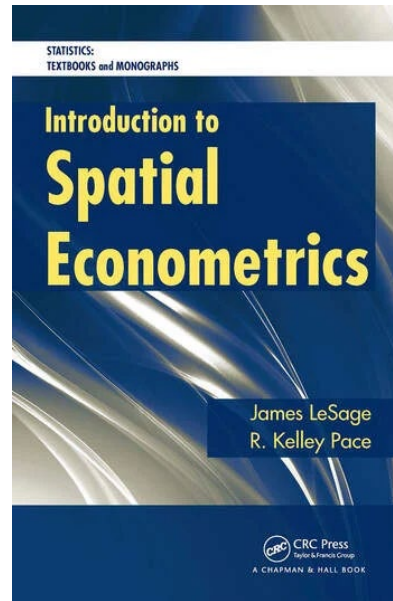
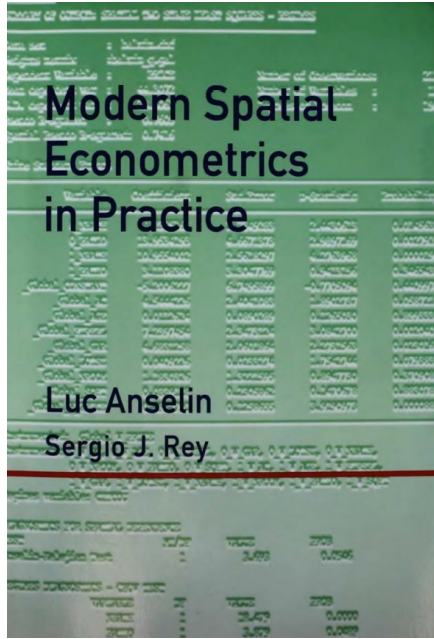
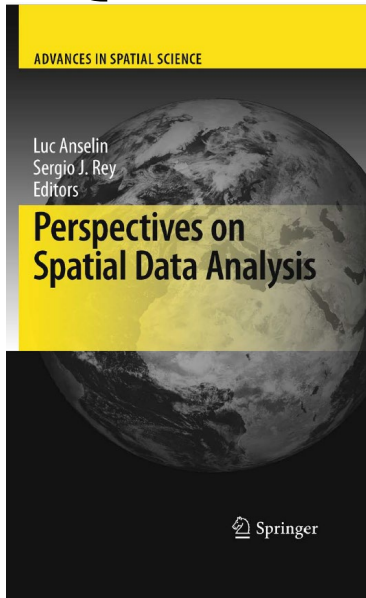
Life Expectancy and its determinants

<http://dagitty.net/dags.html?id=GtvICQ>

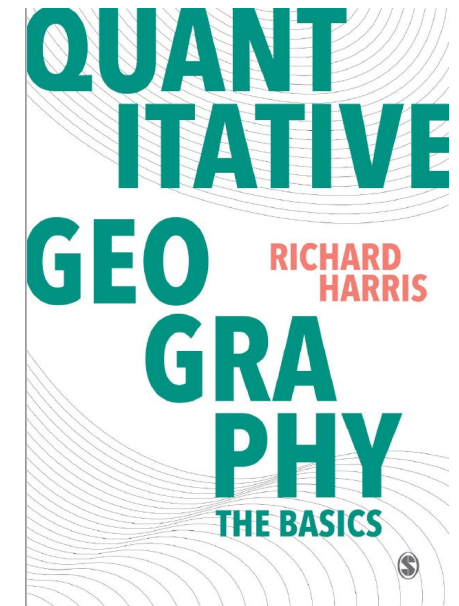
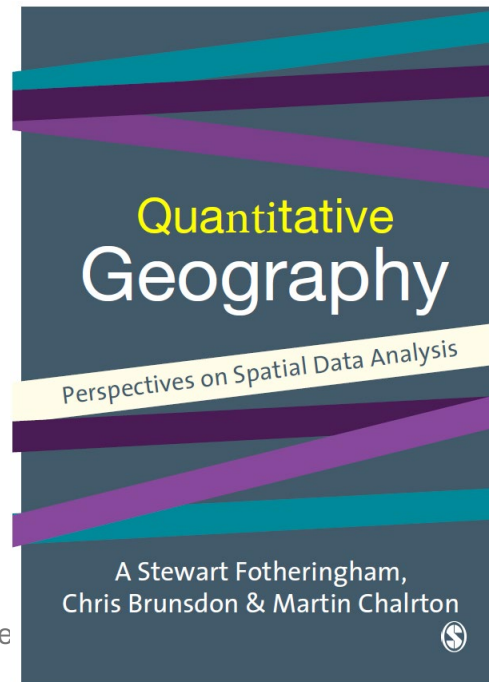
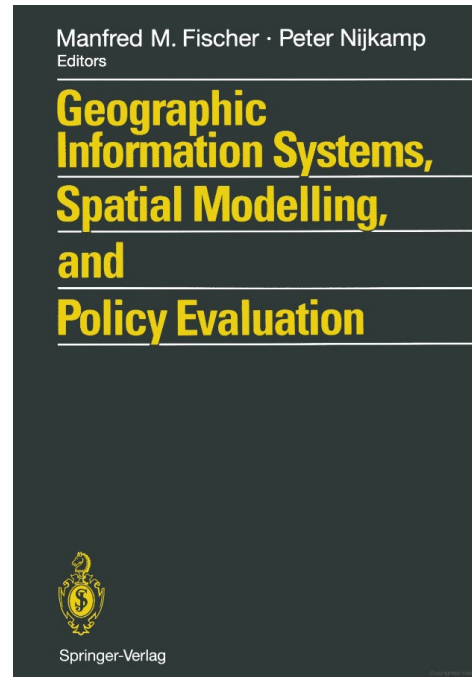


Simplest model' is: everything relates to everything: the 'saturated' model ('reference' in path analytic/SEM lingo)

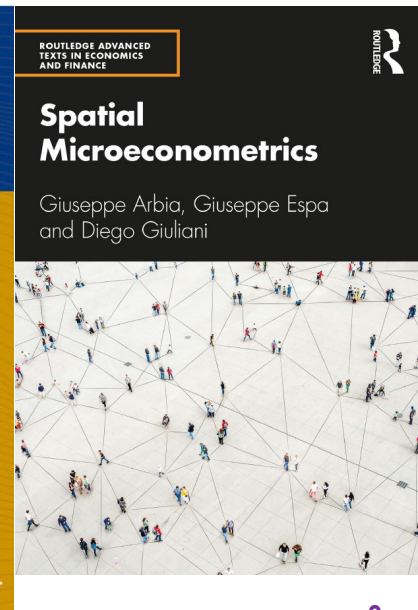
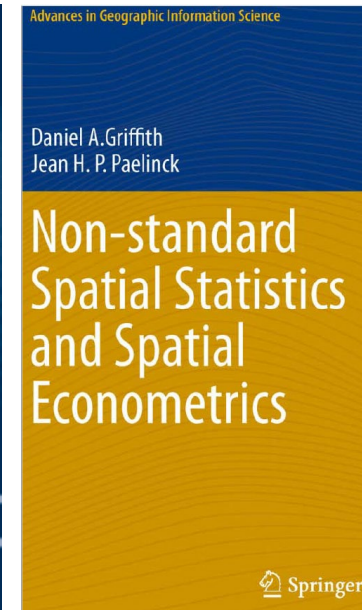
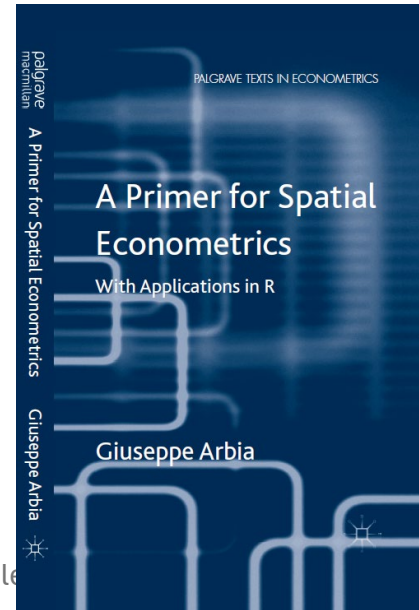
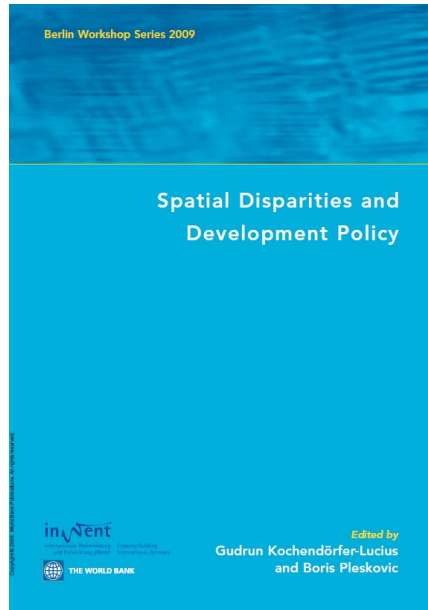
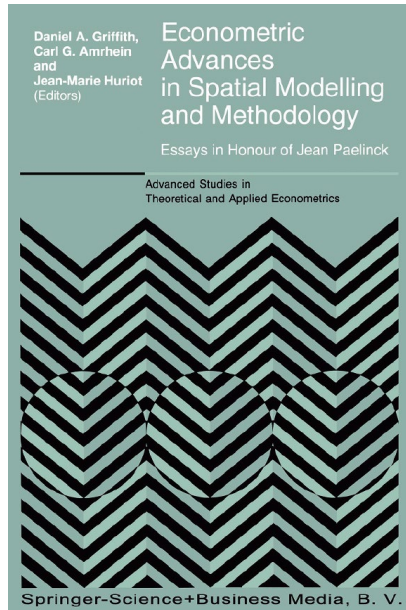
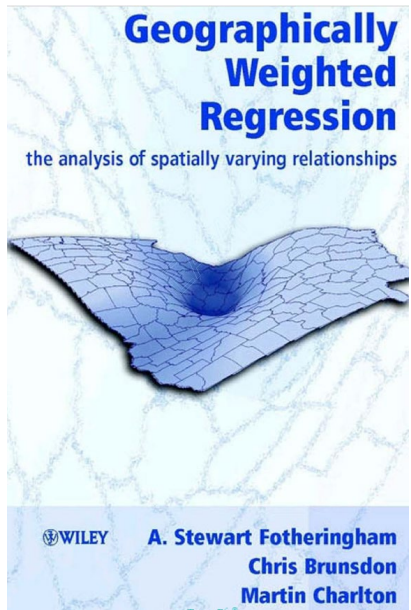
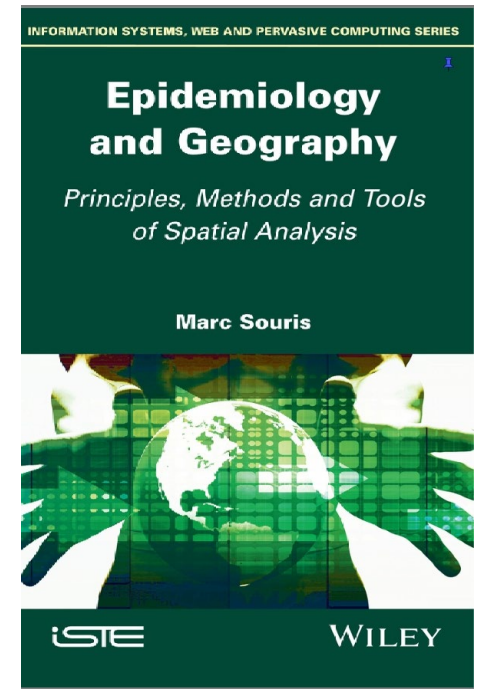
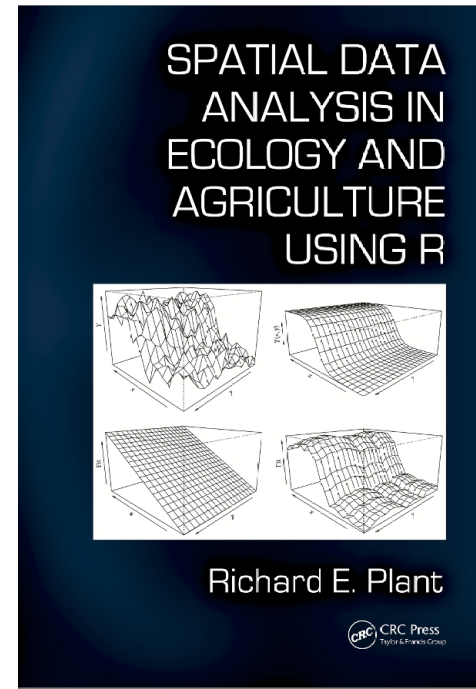
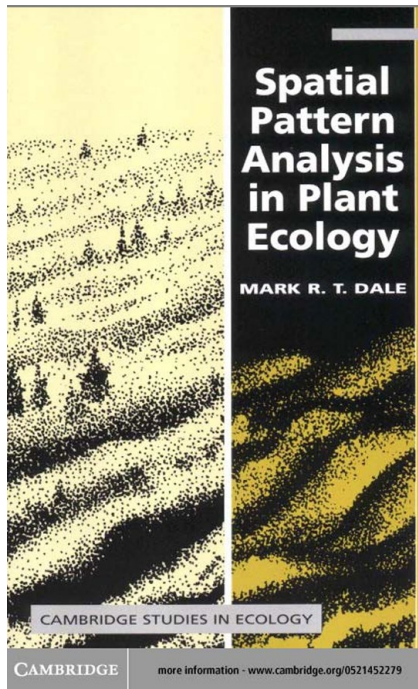
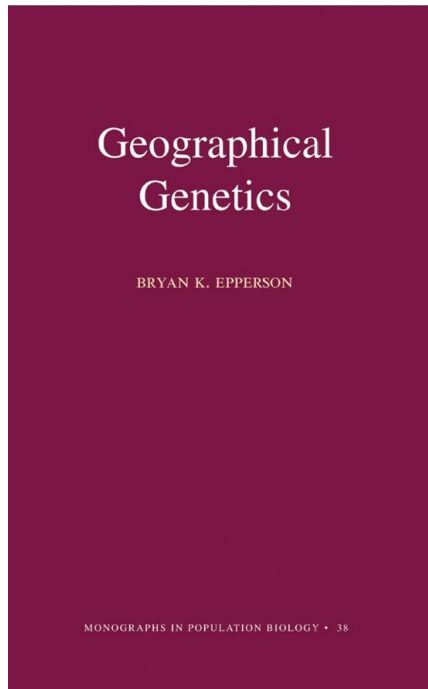
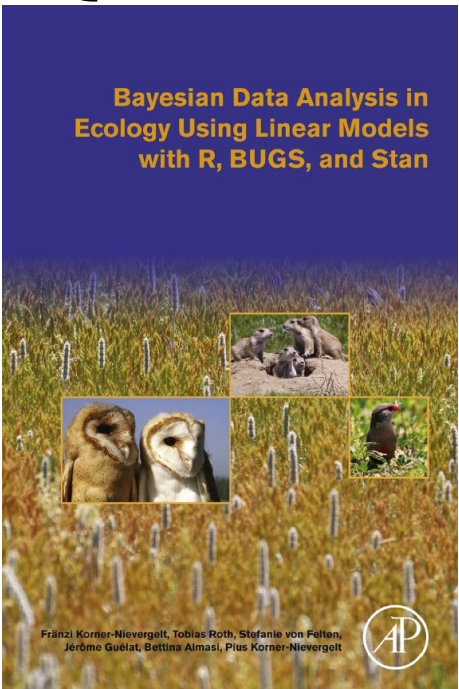
Quick look:



Visual 5 | Colours to print: CMYK | Font(s): Helvetica

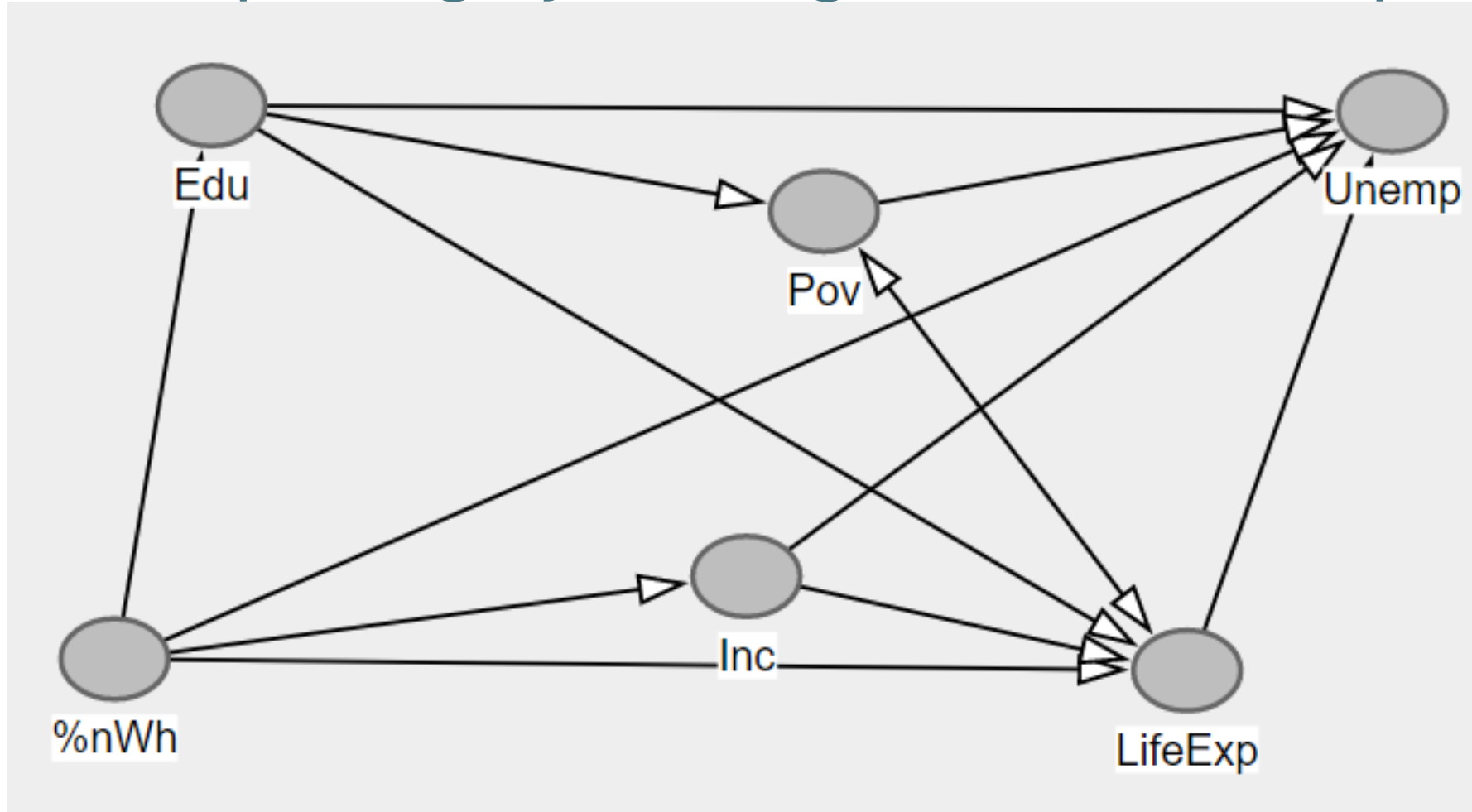


Quick look 2:



Life Expectancy data informed model

<http://dagitty.net/dags.html?id=4TETpl>



There are many factors to consider, of course
http://bit.ly/HD_causal_model, including molecular: [“Scientists Discover a Molecular Switch That Controls Life Expectancy”](#)

Some troubles with spatial/regional/geographic data

A. Averaging to talk about ‘typical region’ does not work :

i. A region with 1 resident with 100y LfEx and another with 100 residents with 80ybLfEx do not yield a 101 aggregate with 90y LfEx.

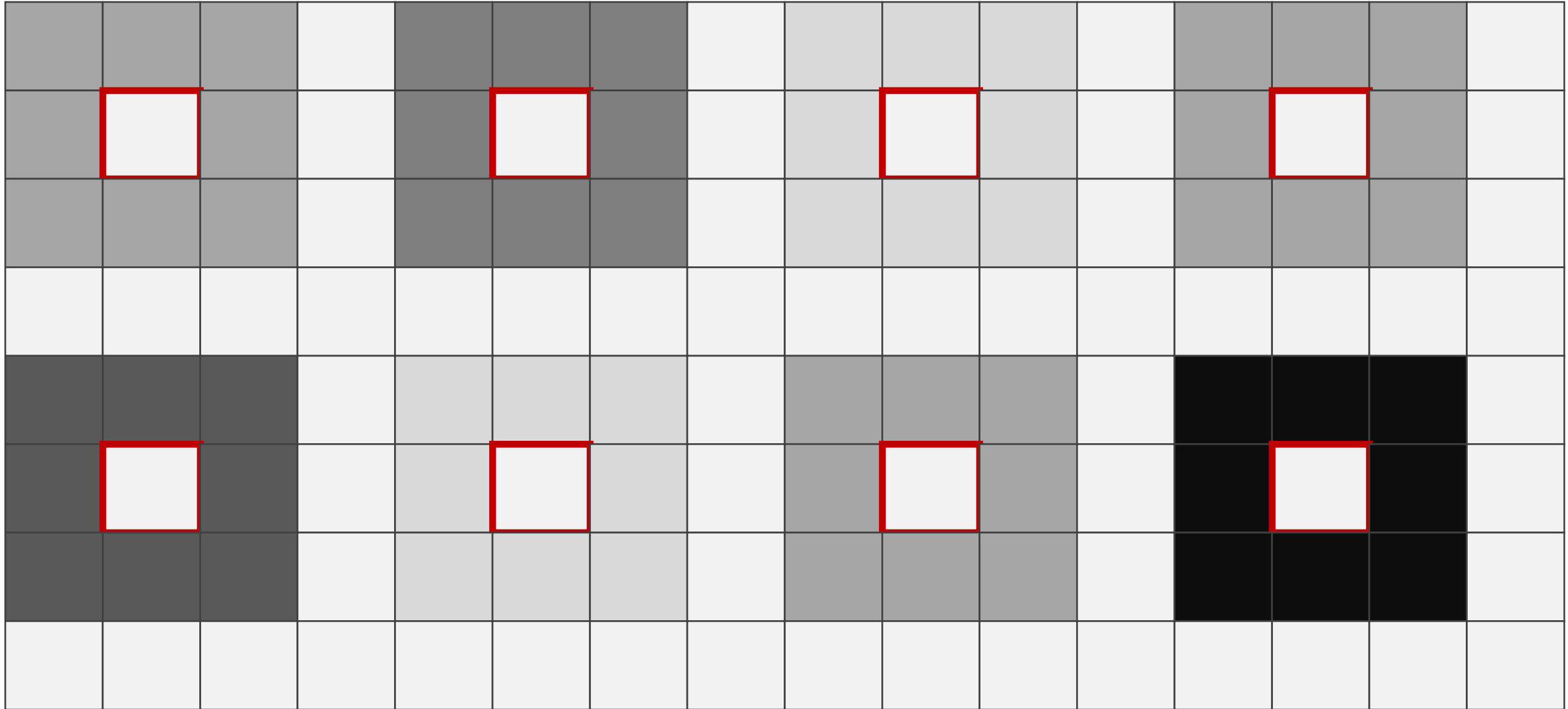
ii. If a region’s LfEx value is *identical* to its neighbors’, then this is too much similarity: much like spousal data, or family data.

B. Clustering within higher level regions due to all-belong-to-higher structure is distinct from clustering due to each-to-its-neighbors spatial structure: there are as many clusters as regions!

* Multilevel modeling does not address the spatial structure, much like it can’t address e.g. friendship relational structure in student-in-classrooms settings.

Intuition for Maximum Moran's I

Haggard, E. A. (1958). [Intraclass correlation and the analysis of variance](#)

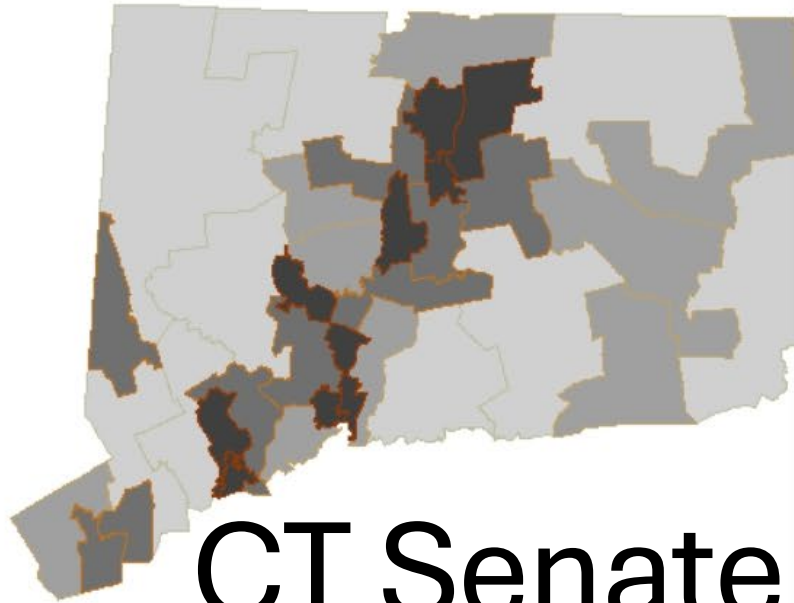


“there is no variation between the scores in any of the [classes [neighbors of red squares]]; rather all the variation is between the [classes [the same #)].”

Quantile: pnw1519sd

- [11.1 : 18.2] (9)
- [19.3 : 27.0] (9)
- [27.0 : 42.7] (9)
- [46.4 : 83.1] (9)

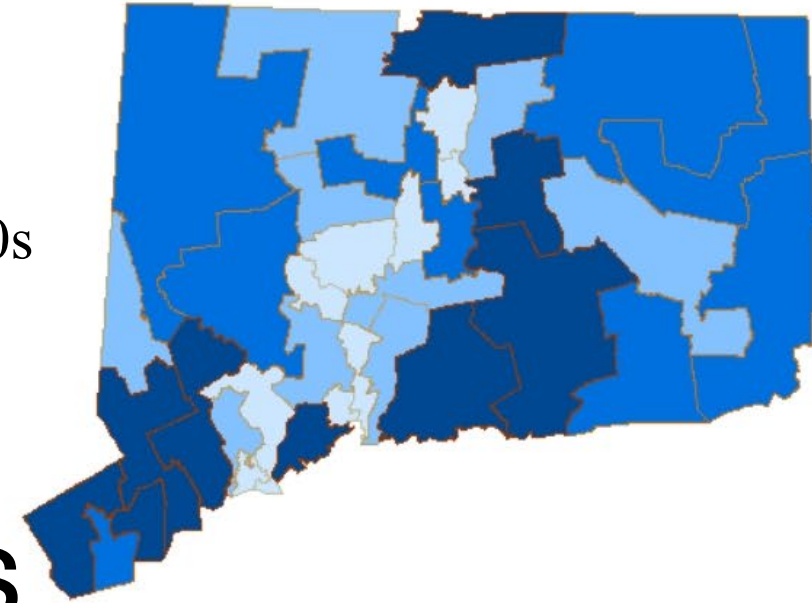
% non-White



Quantile: inc10k1721

- [2.71 : 5.35] (9)
- [5.50 : 6.38] (9)
- [6.47 : 7.78] (9)
- [8.06 : 13.13] (9)

Income \$1,000s

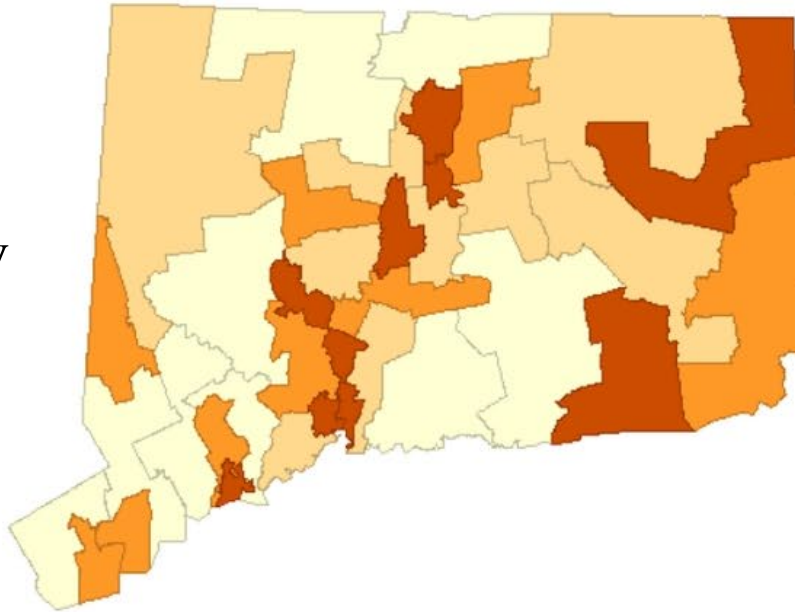


CT Senate Districts

Quantile: ppovsend

- [3.26 : 5.74] (9)
- [5.83 : 8.25] (9)
- [8.78 : 10.21] (9)
- [12.84 : 25.41] (9)

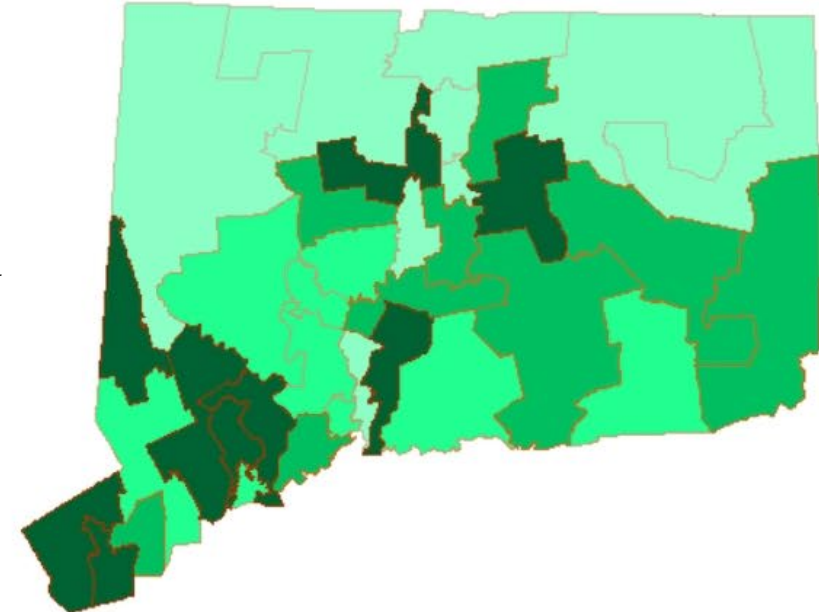
% in Poverty



Quantile: lfexsend

- [69.340 : 73.080] (9)
- [73.386 : 76.594] (9)
- [76.877 : 79.287] (9)
- [79.872 : 82.226] (9)

Life Expectancy

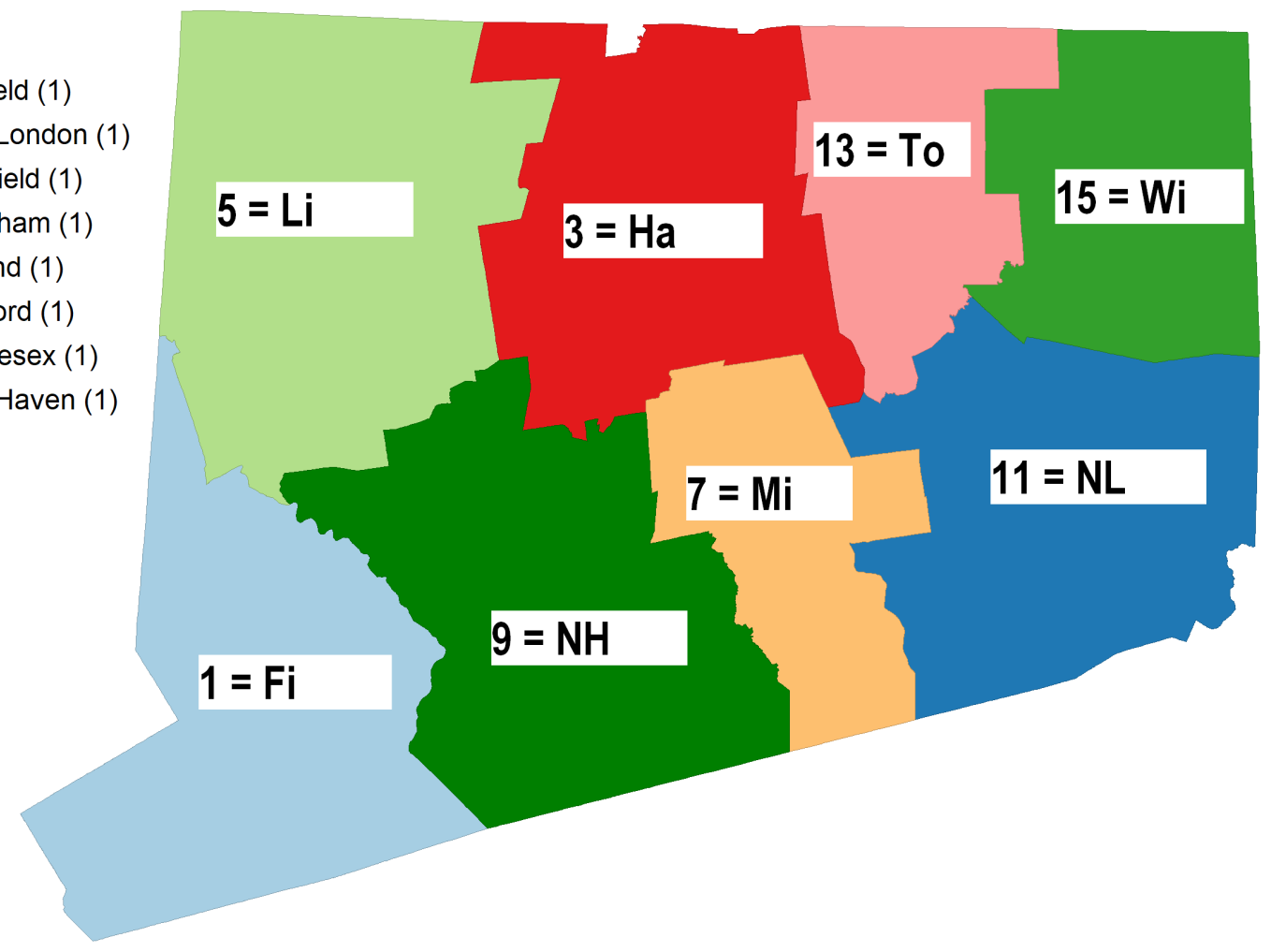


Queen Contiguity Weight Matrix - CT 8 counties

GEOID10

1	2				
5	9				
3	5				
11	5	13	7		
5	3				
1	3	9			
7	3				
11	3	9			
9	4				
3	5	1	7		
11	4				
15	13	3	7		
13	3				
11	15	3			
15	2				
11	13				

- NAME10
- Fairfield (1)
 - New London (1)
 - Litchfield (1)
 - Windham (1)
 - Tolland (1)
 - Hartford (1)
 - Middlesex (1)
 - New Haven (1)



Queen Contiguity Weight Matrix - CT 8 counties

GEOID10
 1 2
 5 9
 3 5
 11 5 13 7 9
 5 3
 1 3 9
 7 3
 11 3 9
 9 4
 3 5 1 7
 11 4
 15 13 3 7
 13 3
 11 15 3
 15 2
 11 13

STANDARDIZE weights:

W_{ij}	9001	9003	9005	9007	9009	9011	9013	9015	
	Fairfield	Hartford	Litchfield	Middlesex	New Haven	New London	Tolland	Windham	Neighbors
Fairfield			0.50		0.50				2
Hartford			0.20	0.20	0.20	0.20	0.20		5
Litchfield	0.33	0.33			0.33				3
Middlesex		0.33			0.33	0.33			3
New Haven	0.25	0.25	0.25	0.25					4
New London		0.25		0.25			0.25	0.25	4
Tolland		0.33				0.33		0.33	3
Windham						0.50	0.50		2 15

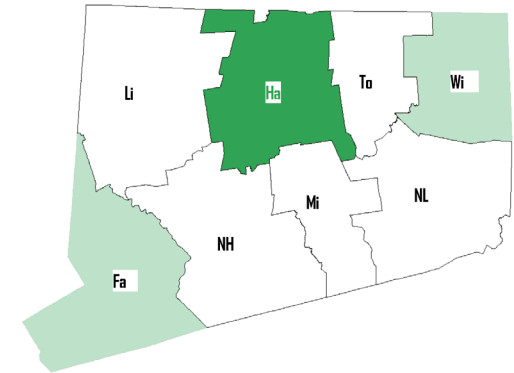
Spelling out the ‘auto’-correlation – CT counties

“In essence, it is a cross-product statistic between a variable and its spatial lag, with the variable expressed in deviations from its mean.” [GeoDa](#)

$$I_Y = \frac{\sum_i \sum_j [(W_{ij} \cdot (y_i - \bar{Y}) \cdot (y_j - \bar{Y})) / S_0]}{[\sum_i (y_i - \bar{Y})^2 / n]}$$

with w_{ij} as the elements of the spatial weights matrix, $S_0 = \sum_i \sum_j w_{ij}$ as the sum of all the weights, and n as the number of observations. For the 8 CT counties, one then would get

$$\begin{aligned} & (Ha [1/5 \cdot (y_{Ha} - \bar{Y}) \cdot (y_{Li} - \bar{Y})] + 1/5 \cdot (y_{Ha} - \bar{Y}) \cdot (y_{Mi} - \bar{Y})] + 1/5 \cdot (y_{Ha} - \bar{Y}) \cdot (y_{NH} - \bar{Y})] + 1/5 \cdot (y_{Ha} - \bar{Y}) \cdot (y_{NL} - \bar{Y})] + \\ & 1/5 \cdot (y_{Ha} - \bar{Y}) \cdot (y_{To} - \bar{Y})] + \\ & [Li \dots] + [Fa \dots] + [NH \dots] + [Mi \dots] + [To \dots] + [NL \dots] + [Wi \dots] +) / 8) / \\ & ((y_{Fa} - \bar{Y})^2 + (y_{Ha} - \bar{Y})^2 + \dots + (y_{Wi} - \bar{Y})^2) / 8) \\ & \text{(if we use the standardized weights, to sum up to 1 per case)} \end{aligned}$$



The ‘clustering’/spatial structure is contained in the Weight Matrix: how the ‘clusters’ are built:

- Each case/region has its own ‘cluster’!
- ‘Clusters’ overlap: same regions can belong to > 1 ‘cluster’!
- There is cyclical influences between ‘members’:

Spelling out the spatial regression – CT counties

A classic regression $Y_i = \alpha. + \beta. \cdot X_i + \varepsilon_i$ would become for spatially connected/nonindependent data e.g., from

$$Y_{\text{Ha}} = \alpha. + \beta. \cdot X_{\text{Ha}} + \varepsilon_{\text{Ha}}, \text{ etc. to:}$$

$$Y_{\text{Ha}} = \rho \cdot (1/5 \cdot Y_{\text{Li}} + 1/5 \cdot Y_{\text{NH}} + 1/5 \cdot Y_{\text{Mi}} + 1/5 \cdot Y_{\text{NL}} + 1/5 \cdot Y_{\text{To}}) + \alpha. + \beta. \cdot X_{\text{Ha}} + \varepsilon_{\text{Ha}}$$

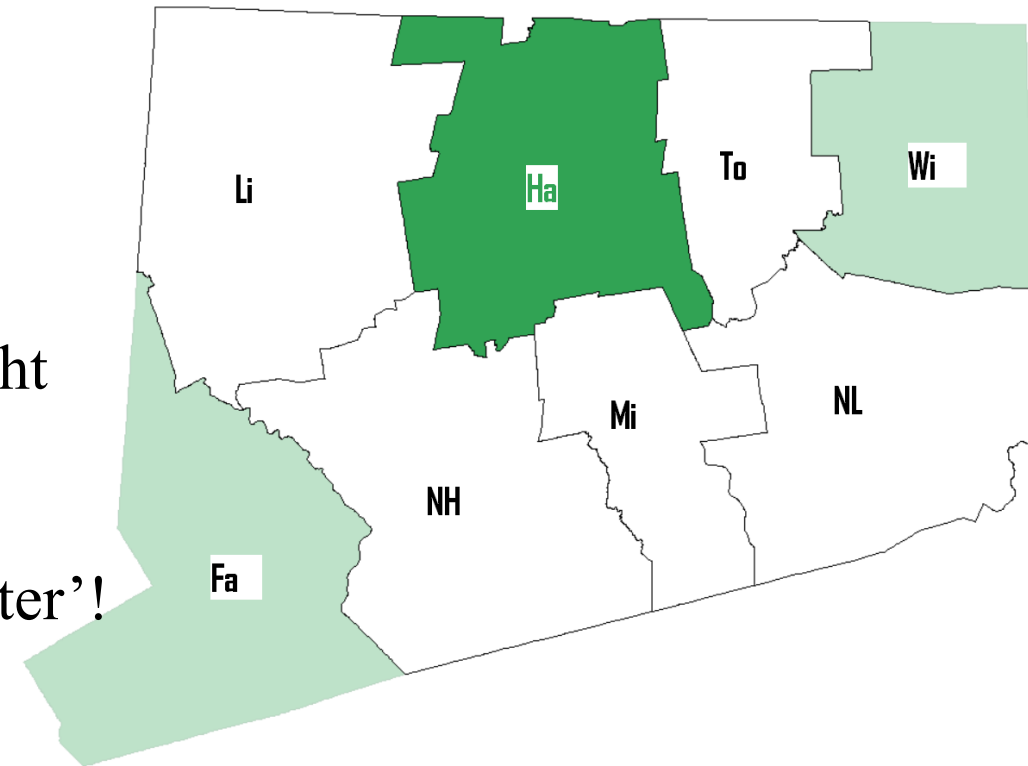
which says that **Ha** has 5 ‘queen’ neighbors,

$$Y_{\text{To}} = \rho \cdot (1/4 \cdot Y_{\text{Ha}} + 1/4 \cdot Y_{\text{Mi}} + 1/4 \cdot Y_{\text{NL}} + 1/4 \cdot Y_{\text{Wi}}) + \alpha. + \beta. \cdot X_{\text{To}} + \varepsilon_{\text{To}},$$

which says that **To** has 4 ‘queen’ neighbors, etc.

The ‘clustering’/spatial structure is contained in the Weight Matrix: how the ‘clusters’ are built:

- Each case/region IS its own ‘cluster’!
- ‘Clusters’ overlap: same regions can belong to > 1 ‘cluster’!
- There is **cyclical influences** between ‘members’:



'Contagion'/interference & causal reasoning

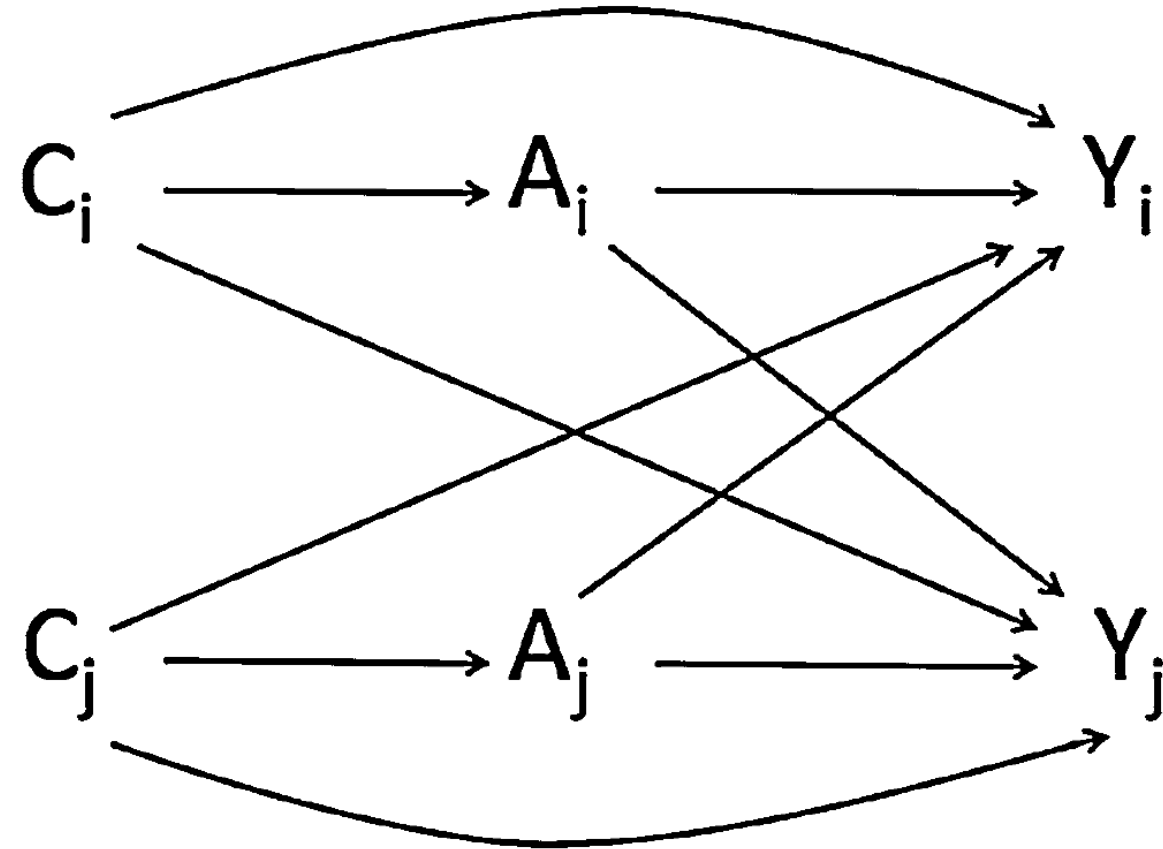
Fig. 5.a

* i and j are 2 'individuals': **regions** here.

* They 'affect' each other ('contagion'): a different type of causal confounding at work.

*** This truly turns patient /clinical/medical health research into public health research.

* The spatial structure adds to this individual DAG (direct acyclic graph) reasoning!

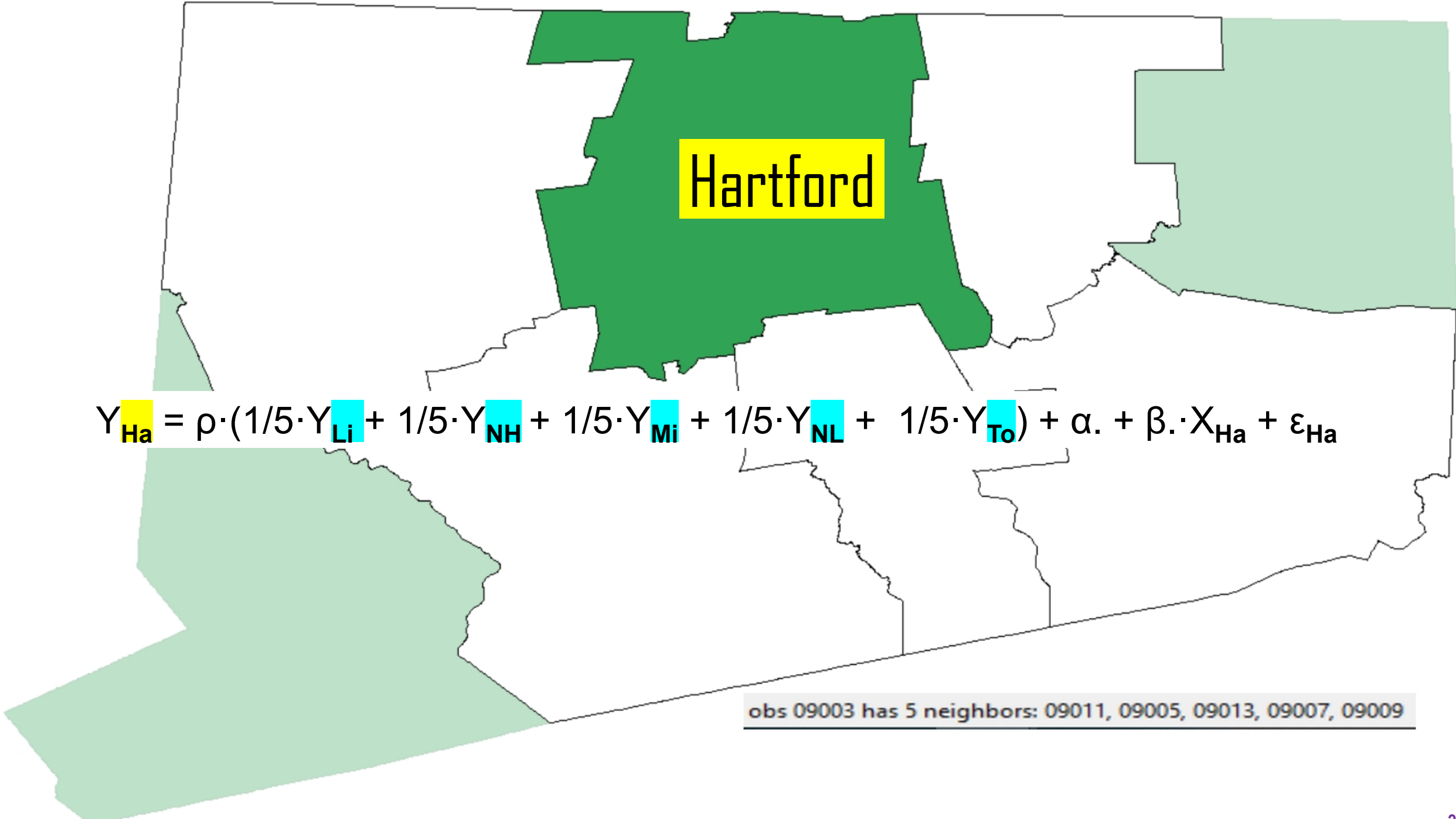




Tolland

$$Y_{\text{To}} = \rho \cdot (1/4 \cdot Y_{\text{Ha}} + 1/4 \cdot Y_{\text{Mi}} + 1/4 \cdot Y_{\text{NL}} + 1/4 \cdot Y_{\text{Wi}}) + \alpha + \beta \cdot X_{\text{To}} + \varepsilon_{\text{To}}$$

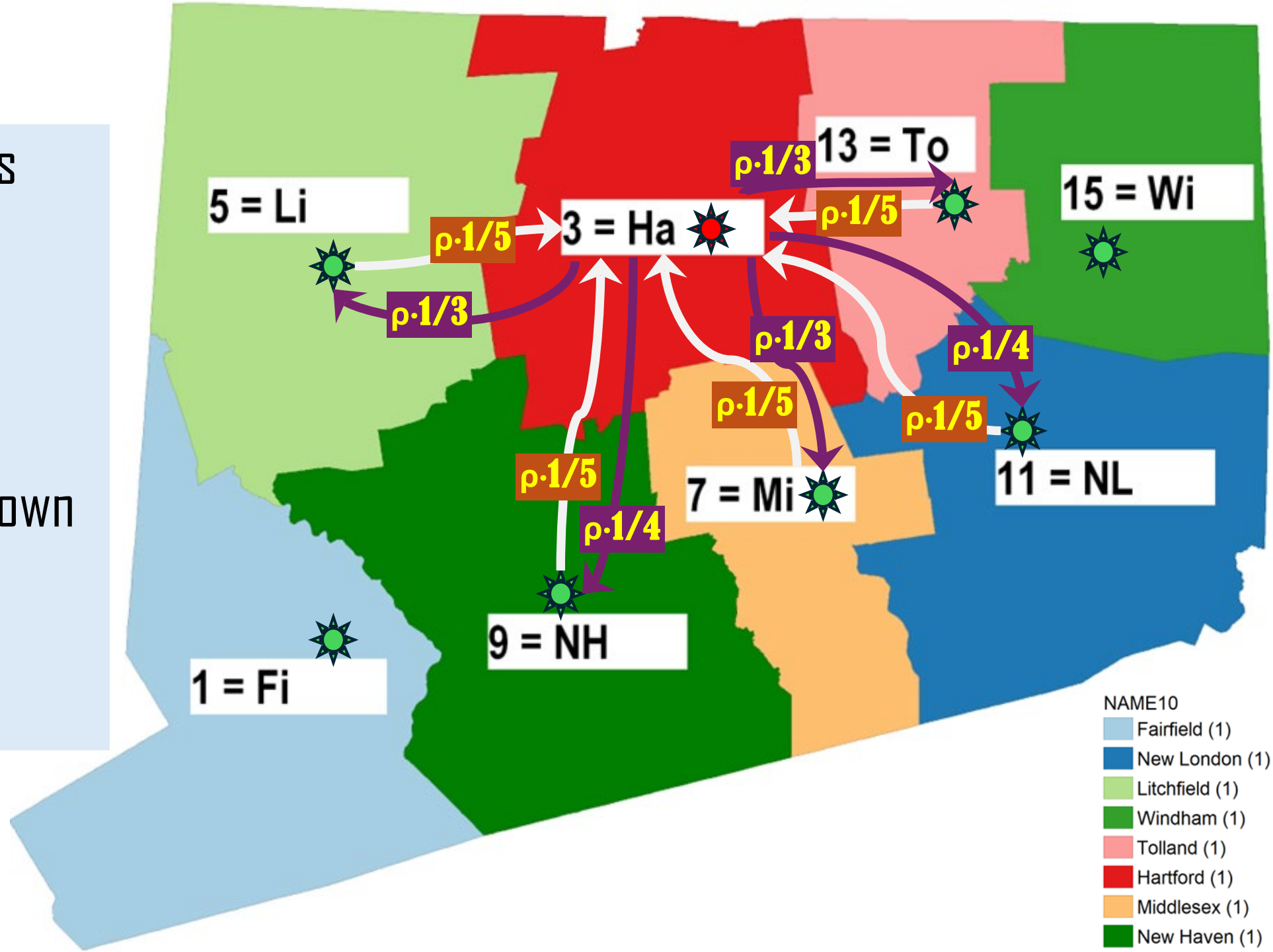
obs 09013 has 3 neighbors: 09011, 09015, 09003



$$Y_{\text{Ha}} = \rho \cdot (1/5 \cdot Y_{\text{Li}} + 1/5 \cdot Y_{\text{NH}} + 1/5 \cdot Y_{\text{Mi}} + 1/5 \cdot Y_{\text{NL}} + 1/5 \cdot Y_{\text{To}}) + \alpha + \beta \cdot X_{\text{Ha}} + \varepsilon_{\text{Ha}}$$

obs 09003 has 5 neighbors: 09011, 09005, 09013, 09007, 09009

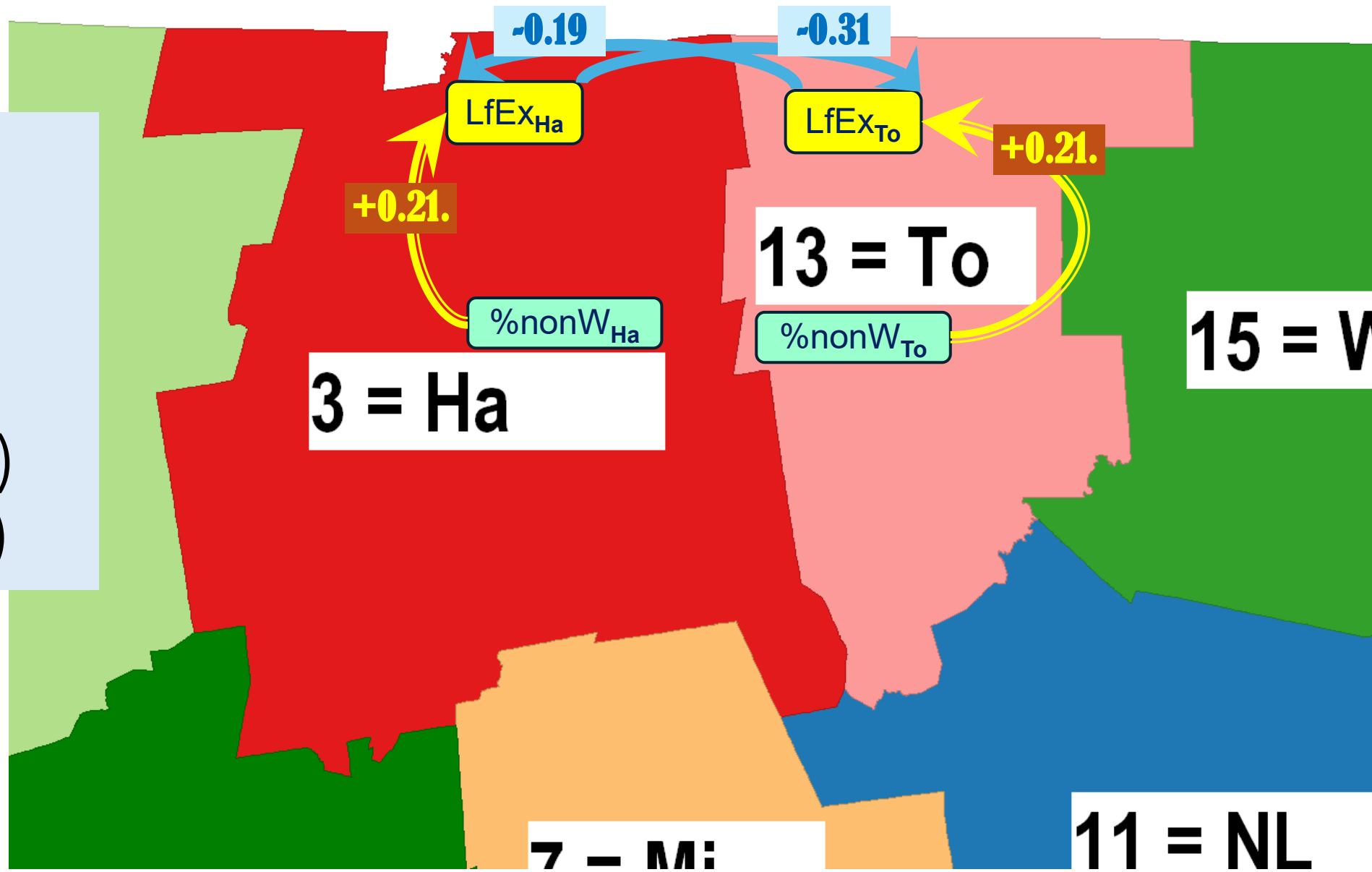
Hartford's LifeExp. is pushed **down** by its neighbors
 $\rho = -0.93$
 $t = -2.522$
But it also pushed down its neighbors, with varying strengths, however.



Spatial 'auto'-
correlation &
causal
reasoning

$$Ha: -0.19 = -0.93 \cdot (1/5)$$

$$To: -0.31 = -0.93 \cdot (1/3)$$



Counties 'individuals': regions here.

Intermediary linguistic clarification

** A correlation is a **same-case** (region, person) & **across/between-variables** statistics:] ‘mutual similarity’ in two sets of numbers: knowing one region’s X tells us something about that region’s Y

** ‘Auto’-correlation is on the other hand **across/between-persons** & **same-variable** statistics: knowing a region’s neighboring regions’ Xs tells us something about its own X

X	Y	
% non-White	Life Exp.	County
11.0	80.3	Litchfield
14.2	81.1	Tolland
14.9	79.2	Windham
15.0	81.3	Middlesex
22.7	79.8	New London
32.5	79.8	New Haven
34.7	82.0	Fairfield
35.2	80.0	Hartford
22.5	80.4	Means

X	Y	
% non-White	Life Exp.	County
14.9	79.2	Windham
22.7	79.8	New London
32.5	79.8	New Haven
35.2	80.0	Hartford
11.0	80.3	Litchfield
14.2	81.1	Tolland
15.0	81.3	Middlesex
34.7	82.0	Fairfield
22.5	80.4	Means

Naïve correlation

Intuition:

When ordering by 1 variable, the other variable's values 'cluster': all high, or all low.

Upper view kind of supports it: 3 of HIGH LifeExp are in LOW %non_White (so we see a 1+3+3+1 pattern in the binary Lo/Hi crosstabulation).

A chi-square test would not find this data pattern statistically significantly different from the null/no relation data pattern (2+2+2+2).

	LifeExp	0	1	
%nWhite	0	1	3	4
	1	3	1	4
		4	4	8

Legend:

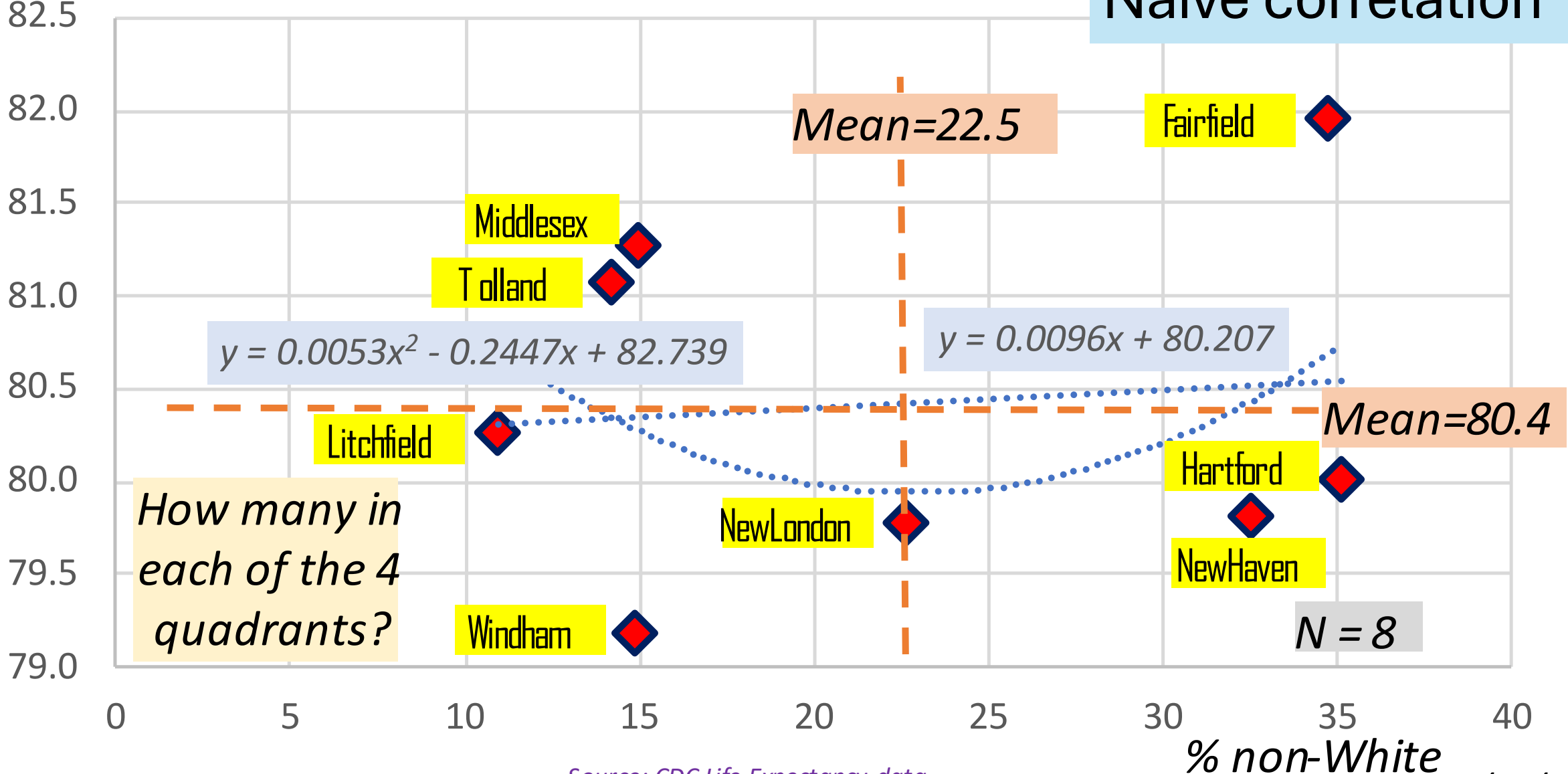
LIGHT color = LOW values

DARK color = HIGH values

Life Expectancy(% non-White residents) in CT, by county

Life Expectancy

Naïve correlation



How many in each of the 4 quadrants?

Source: CDC Life Expectancy data

Life Expectancy(% non-White residents) in CT, by county

Life Expectancy

Naïve correlation

All 3 Green stars PULL DOWN
Tolland's (Red star) LfExp value, AND
pull LEFT %nonWhite value: both l's
<0(NS for n=8): will make appear a p>0

Mean=22.5

Fairfield

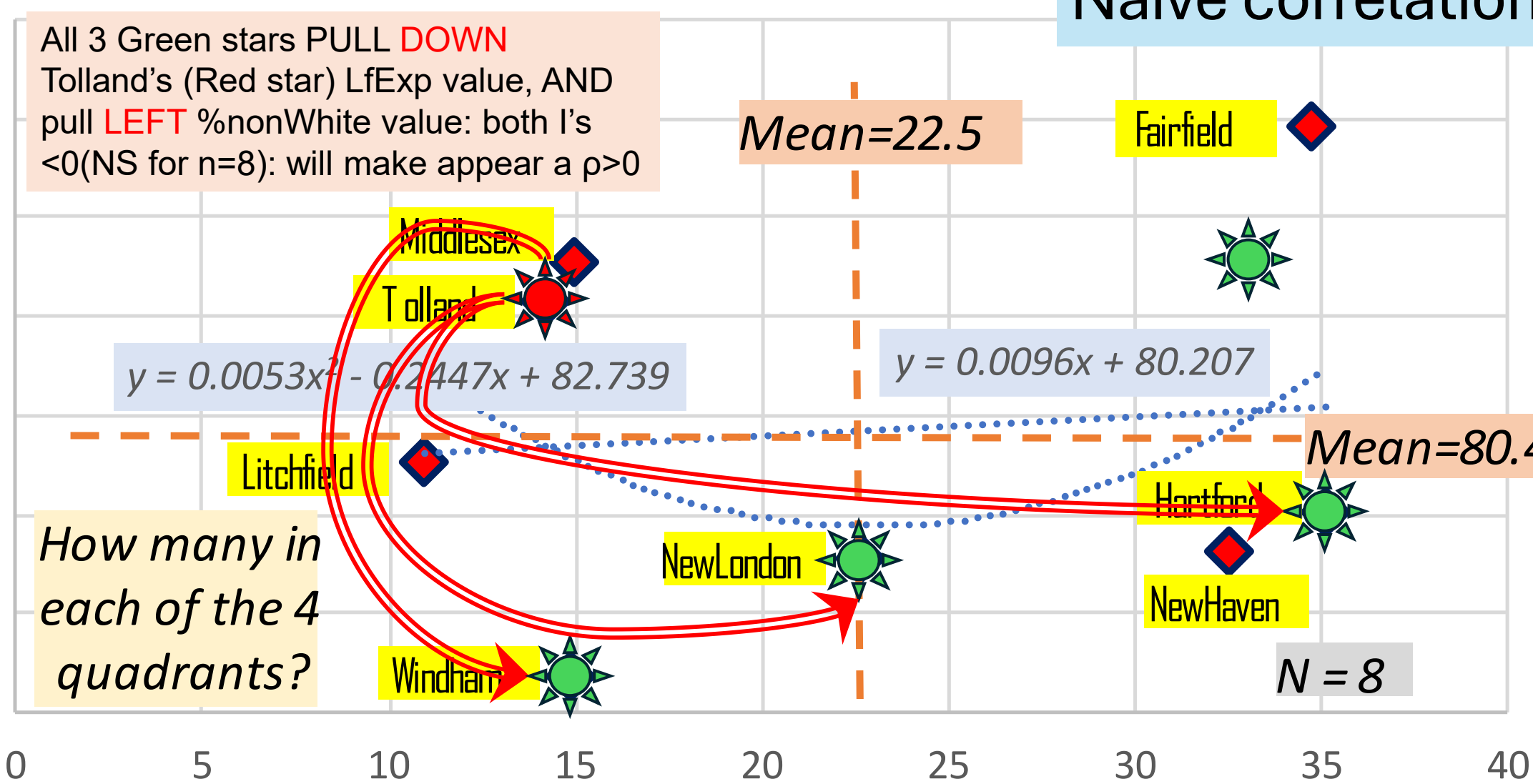
$$y = 0.0053x^2 - 0.2447x + 82.739$$

$$y = 0.0096x + 80.207$$

Mean=80.4

How many in each of the 4 quadrants?

N = 8



Source: CDC Life Expectancy data

CT Small Claims for Medical Debt totals

Year	Numbers of claims	Total filed amounts	Mean amount per docket
2015	10,272	\$ 15,767,136	\$ 1,535
2016	12,056	\$19,382,123	\$ 1,608
2018	12,097	\$20,786,962	\$ 1,718
2019	9,185	\$16,348,638	\$ 1,780

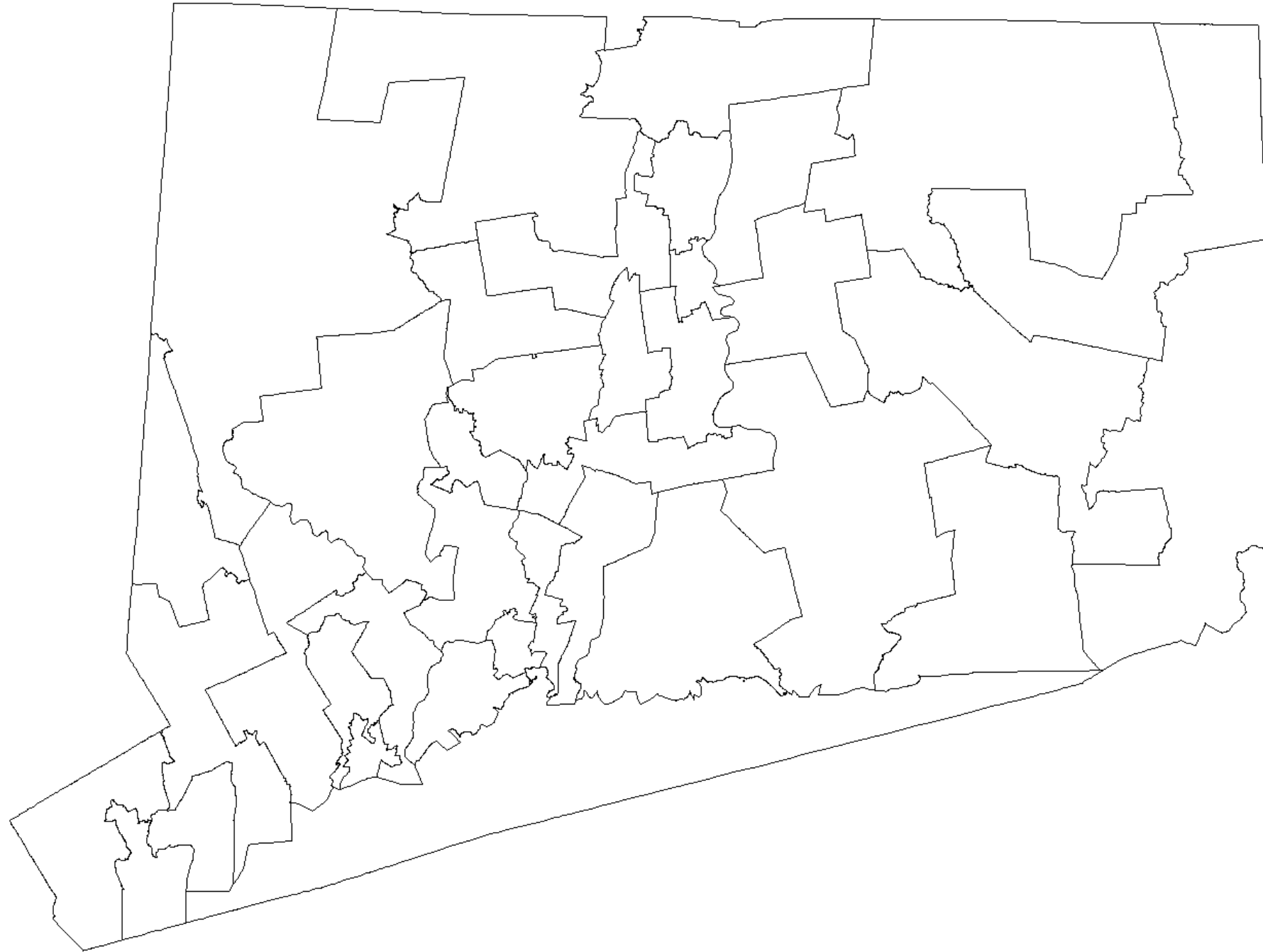
The total number of medical debt small claims in CT, total and average amounts charged per defendant/patient
Notes: The 2017 data did not cover the full year, and is not reported; the claims counts up number of unique 'dockets' or cases filed (multiple family members may appear in the same docket); amounts are shown as 'filed', not as 'awarded' the awarded amounts are 99.2% on the whole from the total amounts filed, in years 1, and 2; years 3 and 4 data did not have amounts awarded.

CT Basic descriptives across 2 geographic/regional layers

	N	Mean
Percent of all people who were nonWhite in 2020 _{CsTr}	820	33.8
CT Senate district	36	32.5
Average annual out of pocket per person on medical care in 2019 _{CsTr}	761	\$1,011
CT Senate district	36	\$ 908
CDC SVI Per capita income estimate, 2014-2018 ACS _{CsTr}	820	\$42,750
CT Senate district	36	\$42,903
Gini Index inequality of household income 2016-2020 _{CsTr}	761	0.427
CT Senate district	36	0.379
Rate of medical debt (per 10,000 residents) in 2019 _{CsTr}	759	28.18
CT Senate district	36	25.75

CsTr = Census Tracts

CT State Senate Districts



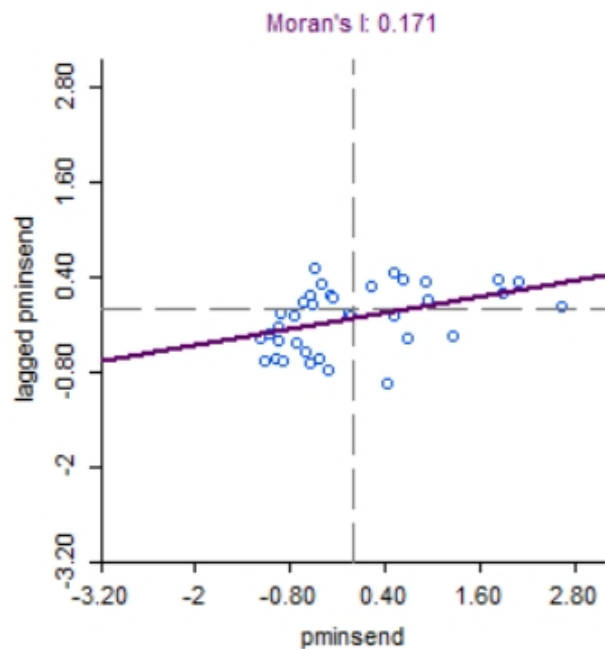
State Senate Districts 'auto'-correlations

$$I_{\text{DebtRate}} = .403, z = 4.186$$

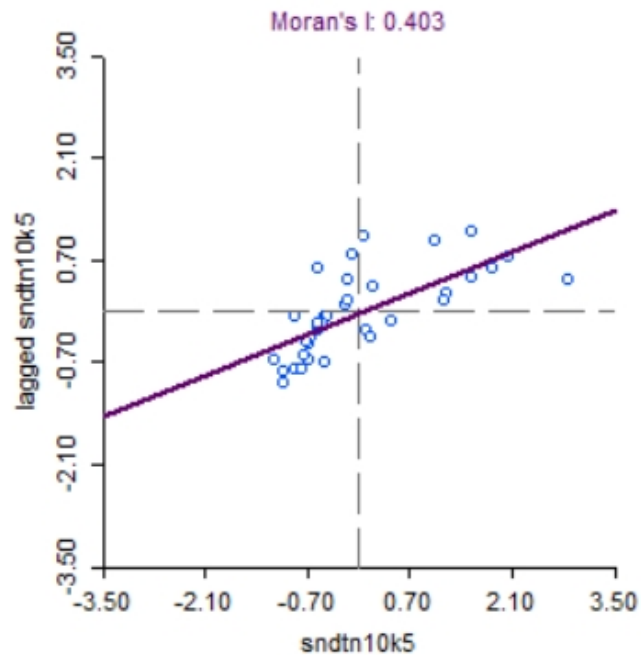
$$I_{\% \text{non-White}} = .171, z = 1.923$$

$$I_{\text{income}} = .329, z = 3.343$$

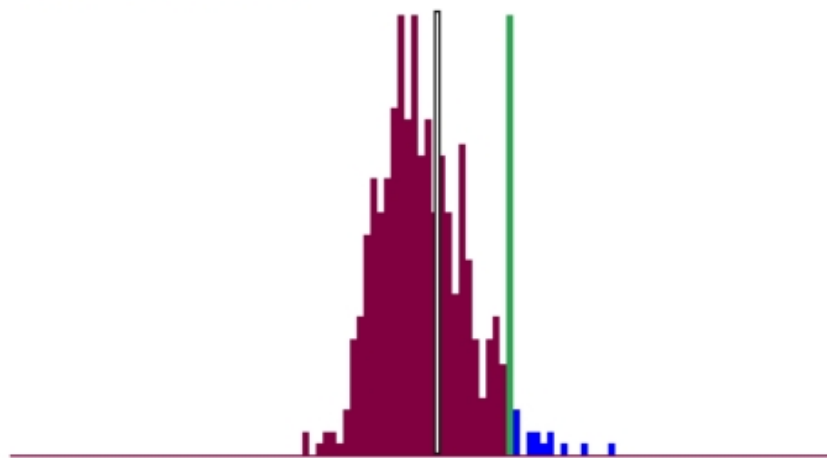
_queen): pminsend



_queen): sndtn10k5

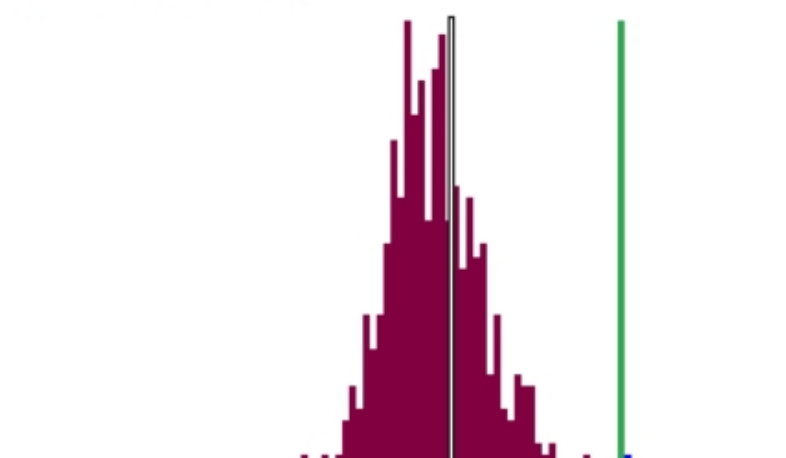


permutations: 499
pseudo p-value: 0.032000



I: 0.1712 E[I]: -0.0286 mean: -0.0305 sd: 0.1049 z-value: 1.9226

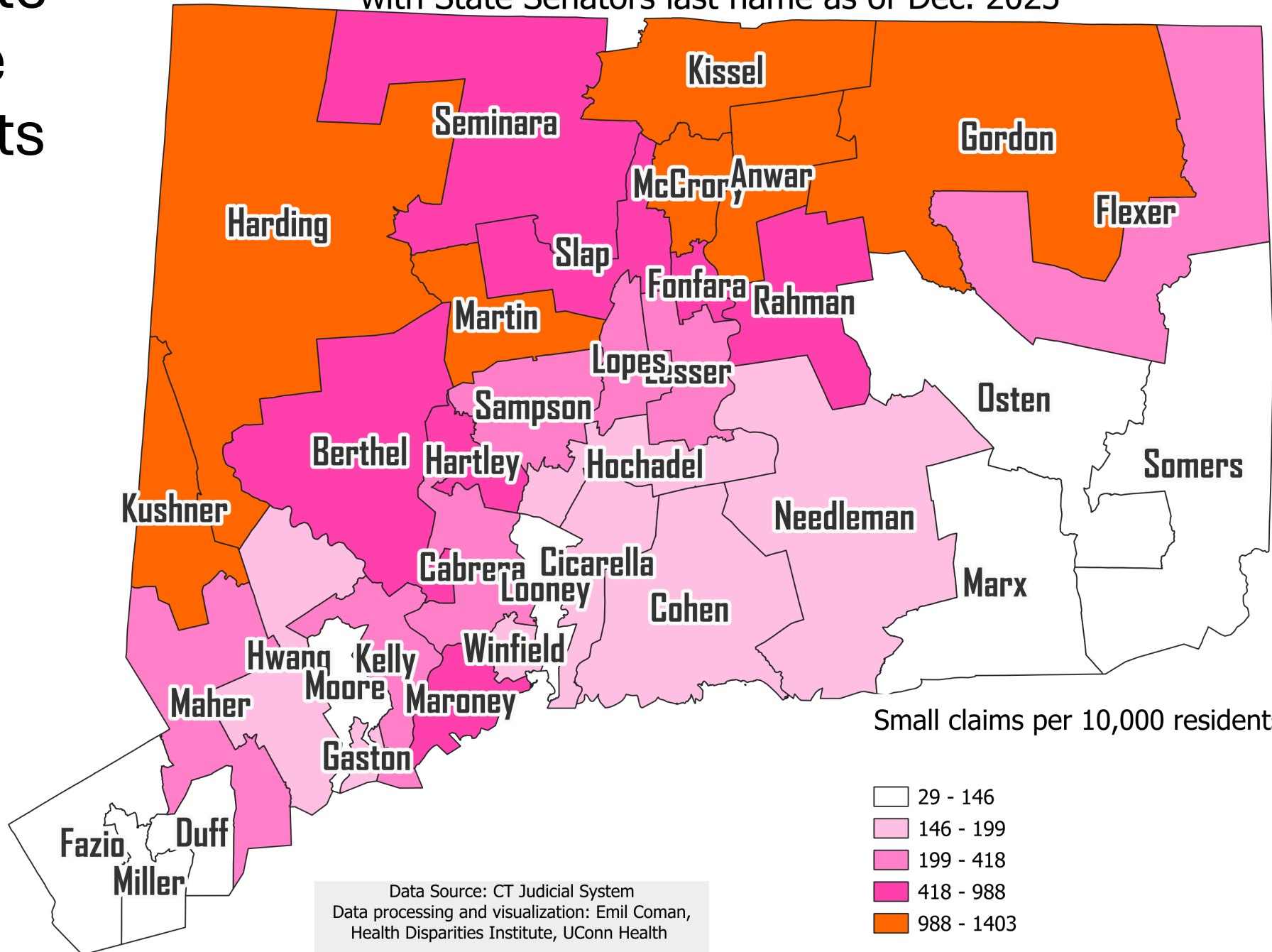
permutations: 499
pseudo p-value: 0.004000



I: 0.4030 E[I]: -0.0286 mean: -0.0334 sd: 0.1043 z-value: 4.1865

CT State Senate Districts

Small Claims medical debt in 2019 in CT by State Senate district, per 10,000 residents, with State Senators last name as of Dec. 2023



Data Source: CT Judicial System
 Data processing and visualization: Emil Coman,
 Health Disparities Institute, UConn Health

2 Predictors of Medical Debt rates in 2019

N = 36, CT state senate districts

	Naïve Beta	Naïve Beta t	Spatial^L Beta	Spatial^L Beta z
% Non-White	0.077	0.432	0.130	1.037
Gini inequality	-91.96	-1.264	108.73 ^{SIG}	-2.120

Notes: ^L - Spatial lag regressions in GeoDa; ^{SIG} - z/t > 1.96.

Seems to suggest that CT state senate districts with more income inequality have a lower debt rate.

Conclusions

1. Estimating effects with spatial data depend require the modeling of spatial 'auto'-correlation, or non-independence.
2. Causal thinking with spatial data forces one to consider two networks: with links between cases (regions), and with links between variables.
3. Spatial data allows for aggregation and mapping of evidence aimed at legislators, or the public.