

## Data Sources and Statistical Results for “Is Decay Contagious? Building Code Violations across New York City”

Jason M. Barr  
August 15, 2025

### Data Sources

*Class C Violations:* New York posts the complete list of building violations over the last decade, downloadable from [https://data.cityofnewyork.us/Housing-Development/Housing-Maintenance-Code-Violations/wvxf-dwi5/about\\_data](https://data.cityofnewyork.us/Housing-Development/Housing-Maintenance-Code-Violations/wvxf-dwi5/about_data). For each 2020 census tract, I added up the total number of violations in 2017-18 and from 2023-24.

*Census Tracts:* Shapefile for NYC downloadable here:

<https://www.nyc.gov/content/planning/pages/resources/datasets/census-tracts>

*Demographics and Income Variables:* For each Census Tract, the data are from the American Community Survey, five-year averages ending in 2023.

- % with Bachelor’s Degree or higher: ACSST5Y2023.S1501
- % Below the Poverty Line: ACSST5Y2023.S1701
- Median Family Income: ACSST5Y2023.S1903
- Race and Ethnicity: ACSDP5Y2023.DP05
- Population: ACSST5Y2023.S1701

*Building Characteristics and Land Use:* Data are from the 2023 PLUTO file, downloadable from <https://www.nyc.gov/content/planning/pages/resources/datasets/mappluto-pluto-change>.

All variables (except borough dummies and distance to the Empire State Building) are calculated at the census tract level. These variables are: Total residential lot area, number of residential units, and average maximum allowable Floor Area Ratio (FAR). Distance to the Empire State Building is as the crow flies from the centroid of each census tract (in degrees).

### Descriptive Statistics

Variable	Obs	Mean	Std. dev.	Min	Max
# Class C Viols, 23-24	2,159	244.6369	372.5788	0	3445
# Class C Viols, 17-18	2,159	90.65262	143.5209	0	1416
BA or +	2,159	38.89402	21.38058	2.9	96.1
Med. Fam. Inc.	2,159	106,177	55,838	18,299	250,000
Population	2,159	3857.925	1899.649	200	16684
Black(%)	2,159	23.25831	26.99946	0	95.6
White(%)	2,159	30.97957	27.77831	0	99.9
Hispanic(%)	2,159	26.88374	21.83048	0	94.6
Asian(%)	2,159	15.35804	17.404	0	91.2
# Res. Units	2,159	1651.252	1203.848	0	20976
Total Lot Area	2,159	2355711	4212569	204745	8.41e+07

MN	2,159	.1315424	.3380709	0	1
BX	2,159	.1523854	.3594773	0	1
BK	2,159	.3483094	.476545	0	1
QN	2,159	.3131079	.4638653	0	1
SI	2,159	.0546549	.2273581	0	1
<hr/>					
Dist. to ESB (deg.)	2,159	.1302244	.0614295	.0028015	.3529803
Avg Allow Max FAR	2,159	1.984724	1.752981	0	10

## Regression Results

Table 1: OLS Regressions. Dep. Var.:  $\ln(1+\text{Violations}$ , 2023-24)

	(1)	(2)	(3)
lnViols17_18	0.920*** (30.19)	0.851*** (26.38)	0.802*** (16.80)
baplus	0.00936 (1.64)	0.00517 (1.41)	
lnMedianInc_fam	-0.376 (-1.78)	-0.220 (-1.60)	
lnPop	0.241 (2.12)	0.0874 (0.89)	
white_pct	-0.00763* (-2.26)	-0.00206 (-0.62)	
black_pct	-0.00814** (-3.47)	-0.00303 (-2.05)	
asian_pct	-0.00462 (-1.56)	-0.000656 (-0.23)	
medIncMwx	-0.168 (-1.89)	-0.285** (-3.36)	
lnResUnits		0.353*** (6.38)	
lnTotalLotArea		-0.191* (-2.57)	
boro2		0.904*** (9.40)	
boro3		0.416** (3.06)	
boro4		0.665** (4.43)	
boro5		0.392 (1.55)	
distESB		-0.577 (-0.63)	
avg_Allow_ResFAR		0.0337 (1.25)	
_cons	1.238*** (14.16)	3.963 (1.82)	3.092 (2.07)
N	2159	2159	2159
R-sq	0.751	0.769	0.792
adj. R-sq	0.751	0.768	0.790
AIC	5555.3	5397.9	5175.6
BIC	5566.6	5420.6	5198.4

t statistics in parentheses. \* p<0.10, \*\* p<0.05, \*\*\* p<0.01. Not variables with zeros are logged using  $\ln(1+ \text{var})$ . Standard errors clustered at the borough level.

Table 2: OLS Regressions. Dep. Var.:  $\ln(1+\text{Violations}, 2023-24)$ , additional specifications

	(1)	(2)	(3)	(4)
lnViols17_18	0.824*** (23.63)	0.804*** (19.86)	0.823*** (25.48)	0.807*** (20.71)
baplus	0.00831** (3.96)	0.00543 (1.89)	0.00405** (3.12)	0.00413 (2.03)
lnMedianInc_fam	-0.331** (-3.59)	-0.226 (-1.75)	-0.343* (-2.73)	-0.227 (-1.59)
poverty_pct	0.000275 (0.12)		-0.00263* (-2.33)	
lnPop	0.385* (2.65)	0.0904 (0.90)	0.0394 (0.49)	0.0886 (0.86)
lnTotalLotArea	-0.145 (-1.88)	-0.208* (-2.52)		-0.195* (-2.48)
boro2	0.905*** (7.34)	0.863*** (9.49)	0.984*** (6.29)	0.921*** (7.77)
boro3	0.315** (3.01)	0.336** (3.58)	0.454* (2.73)	0.401** (3.16)
boro4	0.654** (4.16)	0.587*** (5.32)	0.729** (3.45)	0.686** (4.21)
boro5	0.365 (1.54)	0.332 (1.43)	0.320 (1.31)	0.381 (1.48)
distESB	-1.764 (-1.90)	-0.798 (-0.86)	-1.653 (-1.44)	-0.662 (-0.64)
lnResUnits		0.366*** (5.79)	0.296*** (8.83)	0.348*** (6.29)
white_pct		-0.00178 (-0.97)		
black_pct		-0.00261*** (-7.44)		-0.00191 (-1.89)
medIncMwx		-0.260** (-3.80)		-0.285** (-3.52)
asian_pct			0.00213 (1.34)	
avg_Allow_ResFAR			0.0518 (1.46)	0.0344 (1.26)
_cons	3.731** (4.59)	3.387* (2.55)	2.461 (2.12)	3.185 (2.05)
N	2159	2159	2159	2159
R-sq	0.782	0.791	0.787	0.791
adj. R-sq	0.781	0.790	0.786	0.790
AIC	5271.7	5179.3	5221.5	5177.9
BIC	5294.4	5202.0	5244.2	5200.6

t statistics in parentheses. \* p<0.10, \*\* p<0.05, \*\*\* p<0.01. Note: Since median family income by CT is capped at \$250,000, a dummy variable was included for those CTs at the max as an additional control variable. Standard errors clustered at the borough level.

## A Note on Estimation

I include a spatial error term in the spatial autoregression (SAR) and estimate a spatial error coefficient ( $\lambda$ ). According to Esri, “[SAR] is similar to the ordinary least squares regression formula, in which a dependent variable (y) is predicted by a set of explanatory variables (x) and coefficients ( $\beta$ ). However, the residual term (u) is modeled by a different regression equation. This second regression predicts the residual using a spatial autoregressive parameter  $\lambda$  (lambda) and a spatial weights matrix (W), along with its own residual term ( $\varepsilon$ ).

The lambda parameter quantifies the strength of the spatial dependence in the error term and measures how much one location’s error term influences the error terms of its neighbors. The SEM works by filtering out spatial autocorrelation from each of the variables in the model and performing a regression on the spatially filtered variables. As a result, the coefficient estimates are not as affected by the spatial autocorrelation in each variable.” Note, however, that to be conservative, I don’t include the estimate of  $\lambda$  in the discussion of the “contagion” effect. But given the positive estimate, it suggests that “shocks” to surrounding communities will continue to positively impact the number of Class C violations over time in each CT.

Table 3: SAR Regressions. Dep. Var.:ln(1+Violations, 2023-24)

	(1)	(2)	(3)	(4)	(5)
lnViols17_18	0.777*** (49.05)	0.758*** (31.68)	0.756*** (31.16)	0.752*** (31.85)	0.747*** (31.24)
lnPop	0.0949* (1.80)	0.0929 (1.64)	0.100* (1.79)	0.0863 (1.60)	0.0683 (1.19)
white_one_pct	0.000581 (0.74)	-0.000194 (-0.17)	-0.000619 (-0.51)		
lnResUnits	0.396*** (10.00)	0.392*** (8.21)	0.378*** (8.11)	0.372*** (8.41)	0.411*** (8.57)
boro2	0.823*** (11.34)	0.812*** (9.69)	0.891*** (9.33)	0.836*** (8.71)	0.699*** (8.12)
boro3	0.313*** (5.13)	0.294*** (3.87)	0.339*** (4.27)	0.277*** (3.62)	0.209*** (2.64)
boro4	0.599*** (8.84)	0.544*** (6.06)	0.615*** (6.63)	0.566*** (6.09)	0.443*** (4.74)
boro5	0.313*** (2.95)	0.260* (1.76)	0.369** (2.31)	0.309** (2.00)	0.143 (0.94)
lnTotalLotArea	-0.241*** (-7.93)	-0.234*** (-6.03)	-0.211*** (-5.13)	-0.197*** (-5.00)	-0.219*** (-5.83)
avg_yearbuilt		-0.000157 (-0.92)	-0.000126 (-0.73)		
distESB			-1.003 (-1.47)	-0.977 (-1.46)	
baplus				0.00356* (1.91)	
lnMedianInc_fam				-0.182*** (-2.58)	-0.0910 (-1.38)

poverty_pct			-0.000216 (-0.09)	-0.000350 (-0.14)
medIncMwx			-0.242*** (-2.84)	-0.192** (-2.29)
_cons	0.786* (1.79)	1.149* (1.84)	0.928 (1.44)	2.649*** (2.78)
W				
lnViols23_24	0.0777*** (5.60)	0.0766*** (5.69)	0.0718*** (5.33)	0.0731*** (5.51)
e.lnViols23_24		0.344*** (8.36)	0.350*** (8.53)	0.326*** (7.84)
voils_17_18_ct				0.000219 (0.83)
N	2159	2159	2159	2159

t statistics in parentheses. \* p<0.10, \*\* p<0.05, \*\*\* p<0.01These regressions include spatial right-hand side variables. The weight matrix here is of the “contiguity” variety, where a weight matrix entry is one if it borders a CT, zero otherwise. The weight matrix is normalized so all rows sum to 1. Regression (1) has a spatial AR variable, (2) includes a spatial error term, as do (3) and (4), but with different control variables. Equation (5) also includes a “spatial lag” of building violations in 2017-18 (but are not significant. Standard errors are robust, and the estimation method was Gs2sls.

Table 4: SAR Regressions. Dep. Var.: ln(1+Violations, 2023-24)

	(1)	(2)	(3)	(4)	(5)
lnViols17_18	0.795*** (51.73)	0.782*** (33.32)	0.780*** (32.84)	0.776*** (33.65)	0.772*** (33.18)
lnPop	0.103* (1.93)	0.115** (2.03)	0.124** (2.20)	0.110** (2.07)	0.0875 (1.51)
white_pct	0.000607 (0.76)	-0.000242 (-0.22)	-0.000680 (-0.58)		
lnResUnits	0.399*** (9.99)	0.387*** (7.87)	0.372*** (7.72)	0.360*** (8.09)	0.406*** (8.19)
boro2	0.810*** (10.84)	0.790*** (9.33)	0.879*** (9.07)	0.842*** (8.69)	0.679*** (7.89)
boro3	0.293*** (4.68)	0.255*** (3.35)	0.305*** (3.83)	0.255*** (3.31)	0.154** (1.97)
boro4	0.579*** (8.42)	0.543*** (6.10)	0.623*** (6.70)	0.587*** (6.36)	0.417*** (4.60)
boro5	0.275** (2.56)	0.223 (1.49)	0.344** (2.13)	0.303* (1.94)	0.0921 (0.61)
lnTotalLotArea	-0.243*** (-7.90)	-0.237*** (-6.23)	-0.214*** (-5.33)	-0.198*** (-5.12)	-0.221*** (-5.95)
avg_yearbuilt		-0.000198 (-1.21)	-0.000167 (-1.02)		
distESB			-1.073 (-1.58)	-0.918 (-1.38)	
baplus				0.00520*** (2.78)	

lnMedianInc_fam			-0.225*** (-3.19)	-0.103 (-1.59)
poverty_pct			-0.000275 (-0.11)	-0.000718 (-0.29)
medIncMwx			-0.238*** (-2.75)	-0.186** (-2.19)
_cons	0.800* (1.78)	1.164* (1.91)	0.945 (1.51)	2.981*** (3.17)
W2				
lnViols23_24	0.0529*** (3.73)	0.0559*** (3.70)	0.0502*** (3.31)	0.0557*** (3.74)
e.lnViols23_24		0.438*** (8.76)	0.445*** (8.98)	0.400*** (7.88)
viols_17_18_ct				-0.0000985 (-0.27)
N	2159	2159	2159	2159
2159				
t statistics in parentheses. * p<0.10, ** p<0.05, *** p<0.01. These regressions include spatial right-hand side variables. The weight matrix here is of the "contiguity" variety, where a weight matrix entry is one if it borders a CT or if it borders a CT that borders a CT (i.e., 2-degree contiguity), zero otherwise. The weight matrix is normalized so that all rows sum to 1. Regression (1) has a spatial AR variable, (2) includes a spatial error term, as do (3) and (4), but with different control variables. Equation (5) also includes a "spatial lag" of building violations in 2017-18 (but they are not significant). Standard errors are robust, and the estimation method was Gs2sls.				

Table 5: SAR Regressions. Dep. Var.: ln(1+Violations, 2023-24)

	(1)	(2)
lnViols17_18	0.748*** (30.20)	0.750*** (31.21)
baplus	0.00319* (1.67)	
lnMedianInc_fam	-0.187*** (-2.67)	-0.112* (-1.69)
poverty_pct	-0.0000149 (-0.01)	-0.000329 (-0.13)
lnPop	0.0844 (1.60)	0.0614 (1.11)
avg_Floors	-0.00947 (-0.71)	
lnResUnits	0.349*** (7.93)	0.378*** (8.25)
boro2	0.828*** (7.38)	0.689*** (7.08)
boro3	0.237** (2.41)	0.175* (1.92)
boro4	0.511*** (4.49)	0.371*** (3.68)
boro5	0.256* (1.66)	0.0813 (0.59)

distESB	-1.416** (-2.07)	
lnTotalLotArea	-0.0972** (-2.14)	-0.117*** (-2.58)
medIncMwx	-0.188** (-2.17)	-0.151* (-1.78)
_cons	1.626 (1.57)	1.009 (0.97)
<hr/>		
Widist1		
lnViols23_24	0.117*** (4.62)	0.124*** (3.73)
e.lnViols23_24	0.681*** (9.85)	0.696*** (10.36)
viols_17_18_ct		0.0000796 (0.17)
<hr/>		
N	2159	2159
<hr/>		

t statistics in parentheses. \* p<0.10, \*\* p<0.05, \*\*\* p<0.01. These regressions include spatial right-hand side variables. The weight matrix here is of the "contiguity" variety, where a weight matrix entry is the weighted distance from the centroid ( $w=1/dist$ ) if within a half mile of the CT, zero otherwise. The weight matrix is normalized so that all rows sum to 1. Regression (1) has a spatial AR variable and a spatial error term, (2) includes a "spatial lag" of building violations in 2017-18 (but they are not significant). Standard errors are robust, and the estimation method was Gs2sls.