

The Scene of the Crime Predicting Theft Locations in Chicago

Benjamin Rothschild
MACS 30200
bnroths@uchicago.edu

RESEARCH QUESTION

How accurately can a model predict the locations of thefts, burglaries, and robberies in Chicago?

INTRO

Since January 1, 2016 there have been 128,000 reported crimes, thefts or burglaries in Chicago, a rate of about 500 per day which is 38% of all crimes. Previous research has shown that mathematical models can better predict crime than dedicated crime analysts and by using these models for policing and preventative measures, Los Angeles was able to decrease theft-related crime by 7.4% in a controlled experiment (Courneyea 2003).

In this study I perform a similar analysis in Chicago while using novel datasources that represent community input of neighborhood conditions such as 311 calls, datasets that might indicate high-crime areas such as building violations, red light tickets, liquor licenses and how people describe places through tweets. A system like this could be deployed as a real-time tool to help police departments and cities prioritize their resources.

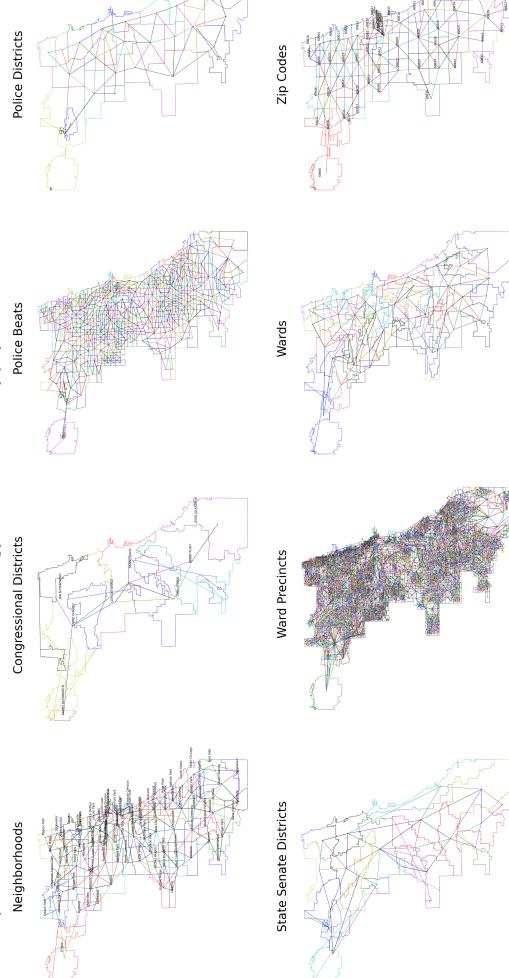
DATA

Table 1: Data Descriptions

	# Obs.	Avg # Obs./day	Source
Crimes - Theft	21,025	253	Chicago Police Department
Crimes - Non-Theft	128,467	324	Chicago Police Department
Building Violations	161,318	55	City of Chicago 311 Requests
311 Garnishment Request	28,537	4	311 Abx Light Out
311 Abx Light On	2,438	0.4	311 Abx Light On
311 Vacant Building	4,433	8	311 Vacant Building
311 Vacant Building Out - No Changes	2,243	4	311 Vacant Building Out - No Changes
Food Inspection - Pass	16,945	33	Chicago Department of Public Health Food Protection Program
Food Inspection - Pass w/Condition	5,674	11	Chicago Department of Public Health Food Protection Program
Food Inspection - Fail	3,342	6	Chicago Department of Public Health Food Protection Program
Liquor Licenses	86,817	964	Illinois Liquor Control Commission
Tweets - Good Sentiment ³	4,941	n/a	Twitter
Tweets - Bad Sentiment ³	1,062	n/a	Twitter
Twitter	183	n/a	Twitter

1. Red light ticket data is only available from 12/2/2013 to 1/3/2014
2. Liquor license data is for current licenses only so it is not shown over time.
3. Twitter dataset was added recently so there is no historical data past what was collected.

Below are maps of all tested boundaries with a line connecting their queen-neighbors. A more detailed explanation of the boundaries and methodology can be found in my paper.



I perform two regressions, an OLS regression which I use as a baseline and a Spatial OLS regression which includes data from spatial neighbors.

In the formulas at left (t_{ij}) denotes the count of an events type in shape j at time t .

To determine the spatial weights I used queen-contiguity weighting matrix.

For my predictions, I used the following cross-validation algorithm to calculate the Mean Squared Error:

1. choose 200 random days within the time-period
2. train model on previous $[t_j, t_i]$ days
3. make prediction for day i

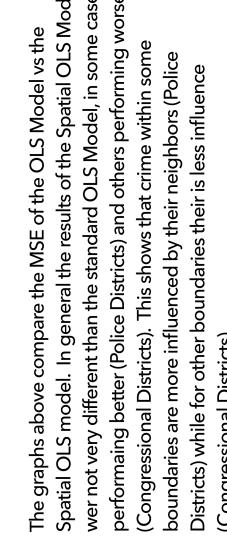
*I tried several values for j and found that in most cases a 28-day "knowledge window" gave a minimum or sufficiently minimum MSE. This methodology is explained more in my paper.

RESULTS & CONCLUSION

Each data point in the table has a timestamp and latitude and longitude attached.

A point-in-polygon calculation is made to place an event into in a shape boundary in the maps below to be used in the regression model.

Some variables such as liquor licenses are included without a time variable in the regression model.



The graphs above compare the MSE of the OLS Model vs the Spatial OLS model. In general the results of the Spatial OLS Model were not very different than the standard OLS Model, in some cases performing better (Police Districts) and others performing worse (Congressional Districts). This shows that crime within some boundaries are more influenced by their neighbors (Police Districts) while for other boundaries their is less influence (Congressional Districts)

Comparing the MSE to the area of the shape used shows the tradeoff between area and MSE when making a prediction.

Normalizing the MSE by the area of the predicting region can show us what "strong" boundaries for crime prediction. This third graph shows that Police Beats, Wards, Neighborhoods and Zip Codes are the strongest boundaries for crime while Ward Precincts and Congressional Districts are weaker boundaries for crime.

$$\begin{aligned} \text{OLS Regression} \\ \text{The } f_{t(i,j)} = \alpha + \beta_1 X_{(t,i)} + \beta_2 X_j + \epsilon_t \\ \text{Spatial OLS Regression} \\ \text{The } f_{t(i,j)} = \alpha + \beta_1 X_{(t,i)} + \beta_2 X_j + \beta_3 \bar{X}_{(t,j \in \text{neighbor}_s)} + \beta_4 \bar{X}_{(t,j \in \text{neighbor}_n)} + \epsilon_t \end{aligned}$$

For my predictions, I used the following cross-validation algorithm to calculate the Mean Squared Error:

1. choose 200 random days within the time-period
2. train model on previous $[t_j, t_i]$ days
3. make prediction for day i

*I tried several values for j and found that in most cases a 28-day "knowledge window" gave a minimum or sufficiently minimum MSE. This methodology is explained more in my paper.