# Time series analysis

## Crime rates and neighborhood demographics

Jason Chang, Jae Hyun, Zhouyang Lian, Sohil Shah, Melanie Tosik

# Business understanding

**Task**
- Predict crime types and locations from spatiotemporal data
- Identify correlations between demographics indicator variables and crime frequencies
- Time series analysis

**Use cases**
- Predict changes in neighborhood "health" based on demographic trends
- Prevent leading causes of mortality, morbidity and social problems among youth in the city
- Improve police dispatch efficiency
- Improve turn-by-turn directions based on type and timing of criminal activity

**Hypotheses**
- Demographic factors are general indicators of a neighborhood's potential susceptibility to crime

# Data understanding

**NYPD Complaint Data Historic**

This dataset includes all valid <u>felony, misdemeanor, and violation crimes</u> reported to the New York City Police Department (NYPD) from 2006 to the end of last year (2016).
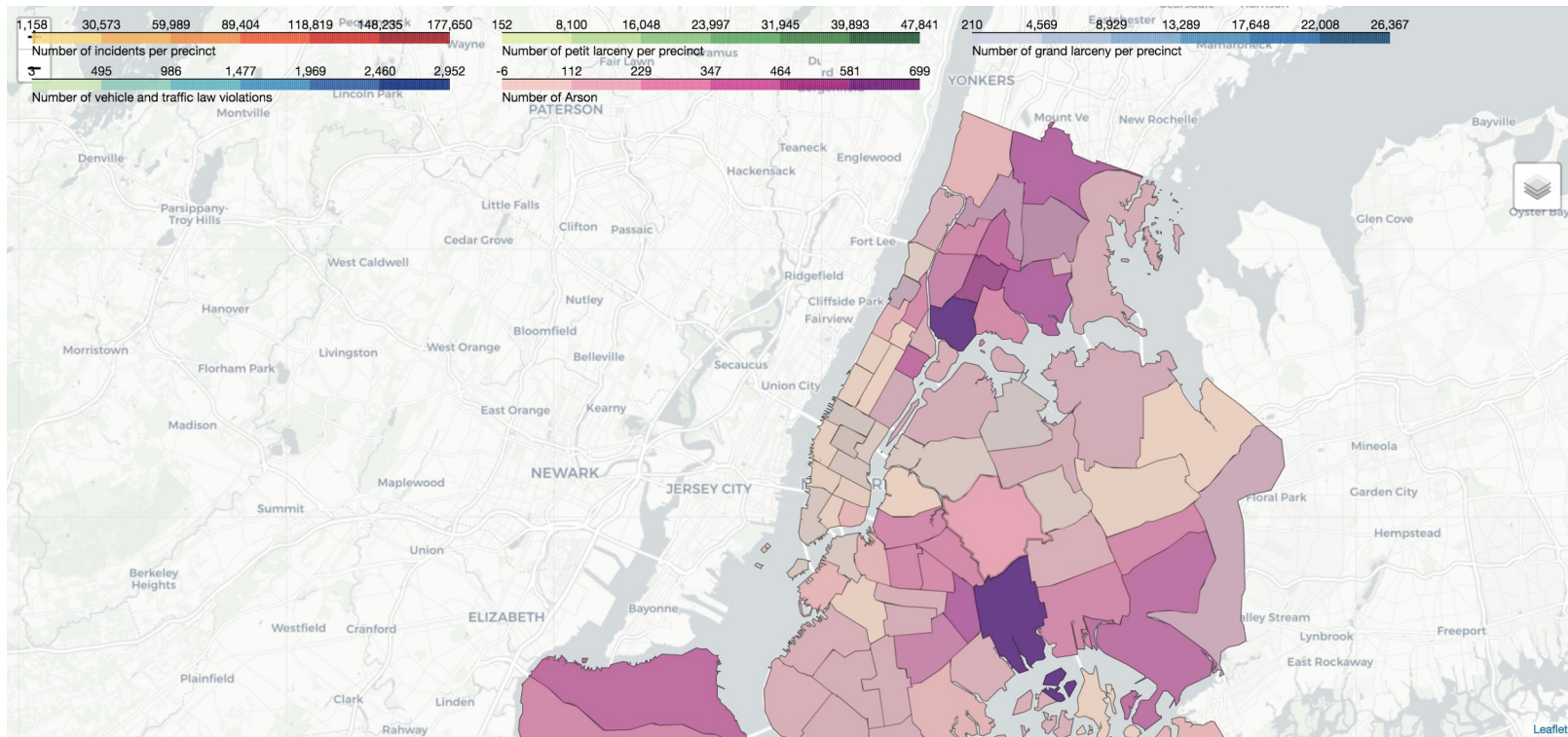
- 2006 to 2016
- 5.58M rows across 24 columns
- Includes date, coordinates, level of offense, description, and responsible jurisdiction

**CoreData.nyc**

New York City's <u>housing and neighborhoods data hub</u>, presented by the NYU Furman Center.
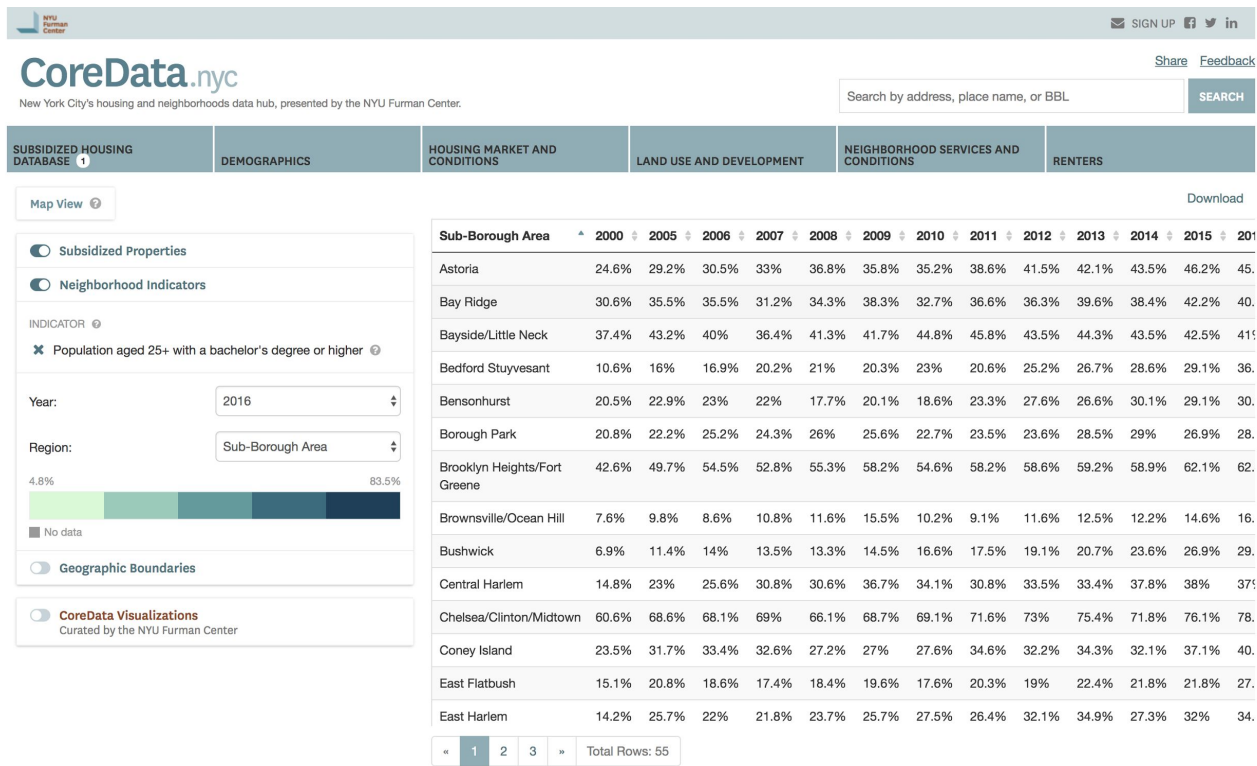
- 20 standardized datasets from city, state, and federal sources
- Demographic data, including household composition, income, education, poverty, and race/ethnicity

# Data understanding



Visualization based on GeoJSON data of police precincts and Leaflet combined with OpenStreetMap

# Data understanding

# Data preparation

**NYPD Complaint Data Historic**

- Feature reduction
    - Missing Values Ratio (MVR)
    - Feature Correlation Threshold (FCT)
- Data aggregation
- Data integration/fusion
- Data transformation
- Visualization

**CoreData.nyc**
- Data normalization
    - Population
    - Crime frequency
- Visualization

# Modeling correlations

**Pearson correlation coefficient**

- Measures linear relationship between two datasets
- Varies between -1 and +1
- -1 or +1 imply exact linear relationship
- 0 implies no correlation
- P-value: probability of uncorrelated system producing data sets with same correlation

**Tools**

- Python
    - pandas
    - sklearn
    - OpenStreetMap
    - Folium (Leaflet.js)
- R
- Tableau

# Evaluation

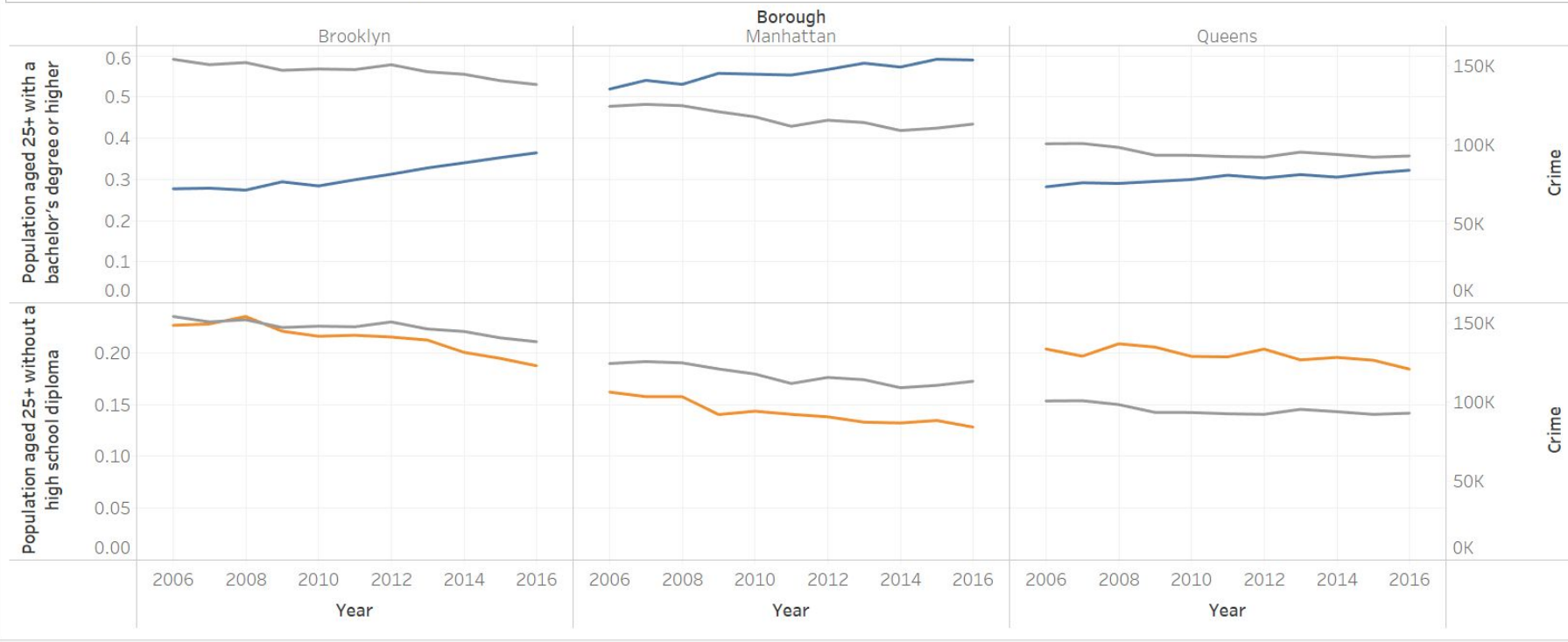**Correlations between crime frequencies and demographic indicators**

- Income diversity ratio                                                [-]*
- Median household income
- Poverty rate                                                          [-]
- Unemployment rate
- **Population aged 25+ with a bachelor's degree or higher**            [-]     Average: -0.795
- **Population aged 25+ without a high school diploma**                 [+]     Average: 0.722
- Population                                                            [-]*
- Population aged 65+                                                   [-]
- Foreign born population
- Racial diversity index

    * 4 of 5 boroughs showed similar trends

# Evaluation



Crime/ education correlations per borough time series

# Introduction to time series analysis

**Definition**

A times series is an ordered sequence of values of a variable at equally spaced time intervals.

**Objectives**

1. Compact description of data,     e.g. classical decomposition: $y_t = T_t + C_t + S_t + I_t$
2. Interpretation,     e.g. long term trends or seasonal variation
3. Forecasting,     e.g. "Twitter predicts the stock market" (Zeng et al., 2010)
4. Control,     e.g. impact of monetary policy on unemployment
5. Hypothesis testing,     e.g. global warming
6. Simulation,     e.g. estimate probability of catastrophic events

**Industry applications**

Economic forecasting, sales forecasting, budgetary analysis, stock market analysis, yield projections, process and quality control, inventory studies, workload projections, or census analysis

# Time series techniques

**Decomposition**
Decompose time series into long term trends, seasonal variation, repeated but non-periodic (cyclic) fluctuations, and residuals (irregular components)

**Forecasting**
Forecast future events based on historic data, based on <u>autoregressive moving average</u> (ARMA) or <u>autoregressive integrated moving average</u> (ARIMA) models

**Clustering**
Partition time series data into groups based on similarity or distance (Euclidean, Manhattan, Hamming, or <u>dynamic time warping</u> (DTW) distance)

**Classification**
Build a classifier based on labeled time series, use it to predict label of unlabeled time series (e.g. k-NN)

# Forecasting and statistical stationarity

**Definition**

A stationary time series is one whose <u>statistical properties are constant</u> over time.

- Stationarity is <u>often required</u> and can be achieved through <u>mathematical transformations</u>
- Stationarized series are <u>easy to predict</u>
- Predictions can be "untransformed" to obtain predictions for the original series
- Stationarizing through <u>differencing</u> is an important part of fitting an <u>ARIMA model</u>

**First difference**

The first difference of a time series is the series of changes from one period to the next.

<u>Example:</u>

If $Y_t$ denotes the value of time series $Y$ at period $t$, first difference of $Y$ at period $t$ is equal to $Y_t$-$Y_{t-1}$.

# ARIMA models

**Definition**

General models for forecasting a time series that can be stationarized through differencing.

- Can be viewed as a "filter" that tries to <u>separate the signal from the noise</u>
- Signal is then <u>extrapolated</u> into the future to obtain forecasts

**Equation**

Linear equation in which the predictors consist of <u>lags of the dependent variable</u> and/or <u>forecast errors</u>.

- Predicted value of $Y$ = constant and/or weighted sum of one or more recent values of $Y$ and/or weighted sum of one or more recent values of the errors
- If predictors consist only of lagged values of $Y$, it is a pure autoregressive model (i.e. a special case of a regression model)

# ARIMA models

**Acronym**

ARIMA stands for **A**uto**R**egressive **I**ntegrated **M**oving **A**verage.
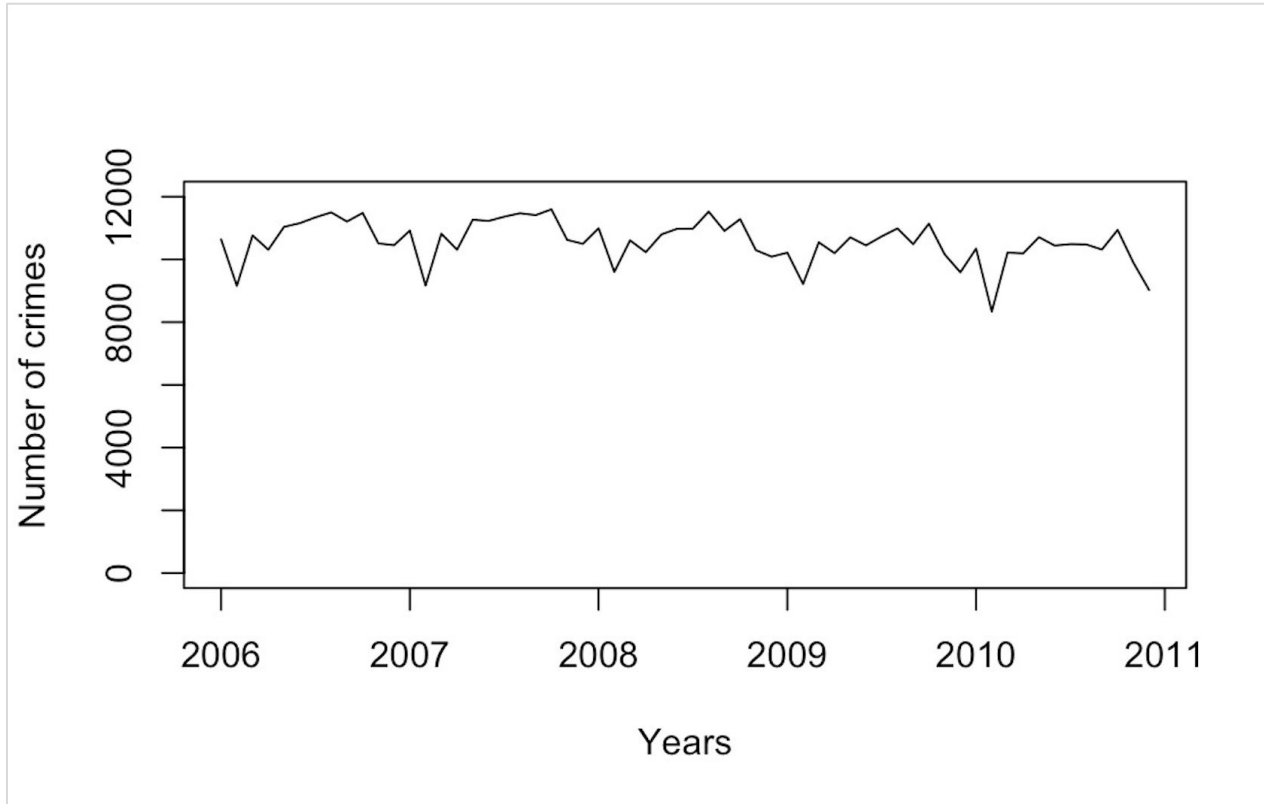
- <u>Lags of the stationarized series</u> in the forecasting equation are called "autoregressive" terms
- <u>Lags of the forecast errors</u> are called "moving average" terms
- Time series which needs to be <u>differenced to be made stationary</u> is said to be an "integrated" version of a stationary series
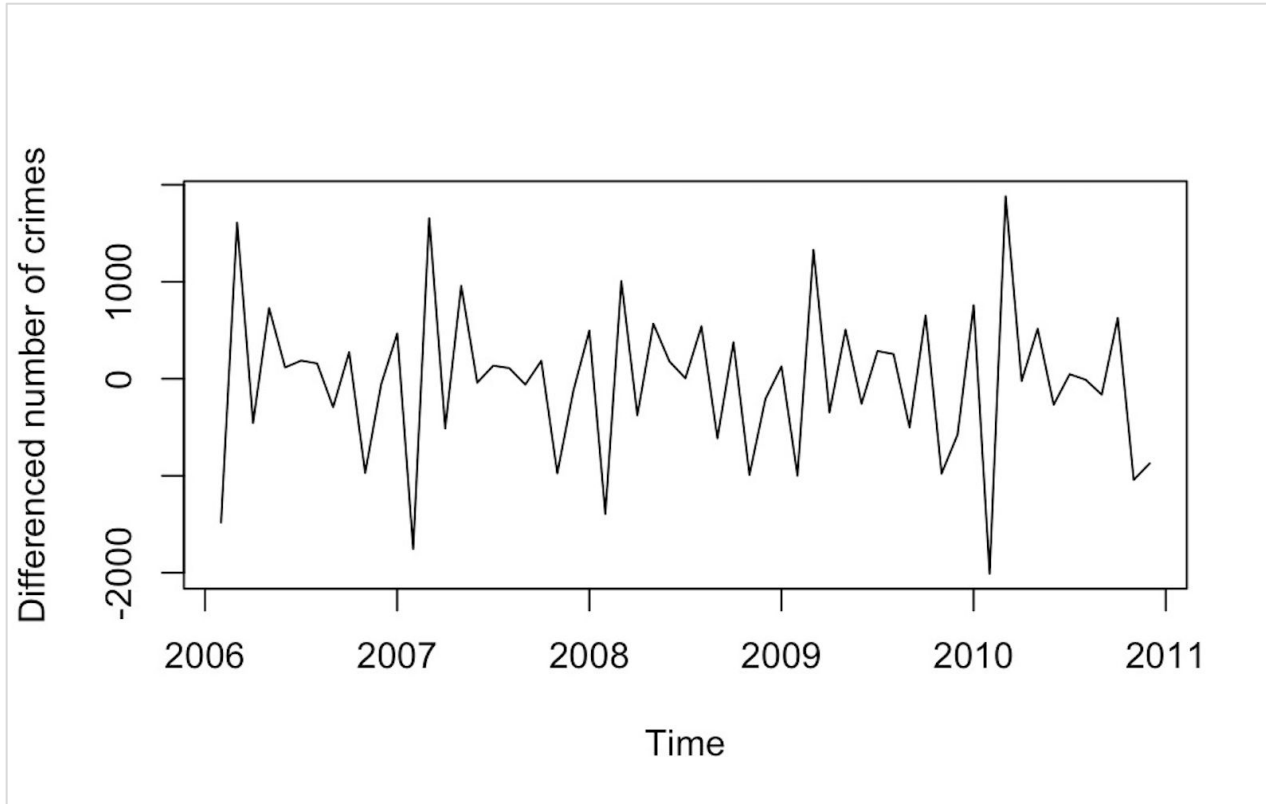
**Equation**

A nonseasonal ARIMA model is classified as an "ARIMA($p,d,q$)" model, where:

- $p$ is the number of autoregressive terms,
- $d$ is the number of nonseasonal differences needed for stationarity, and
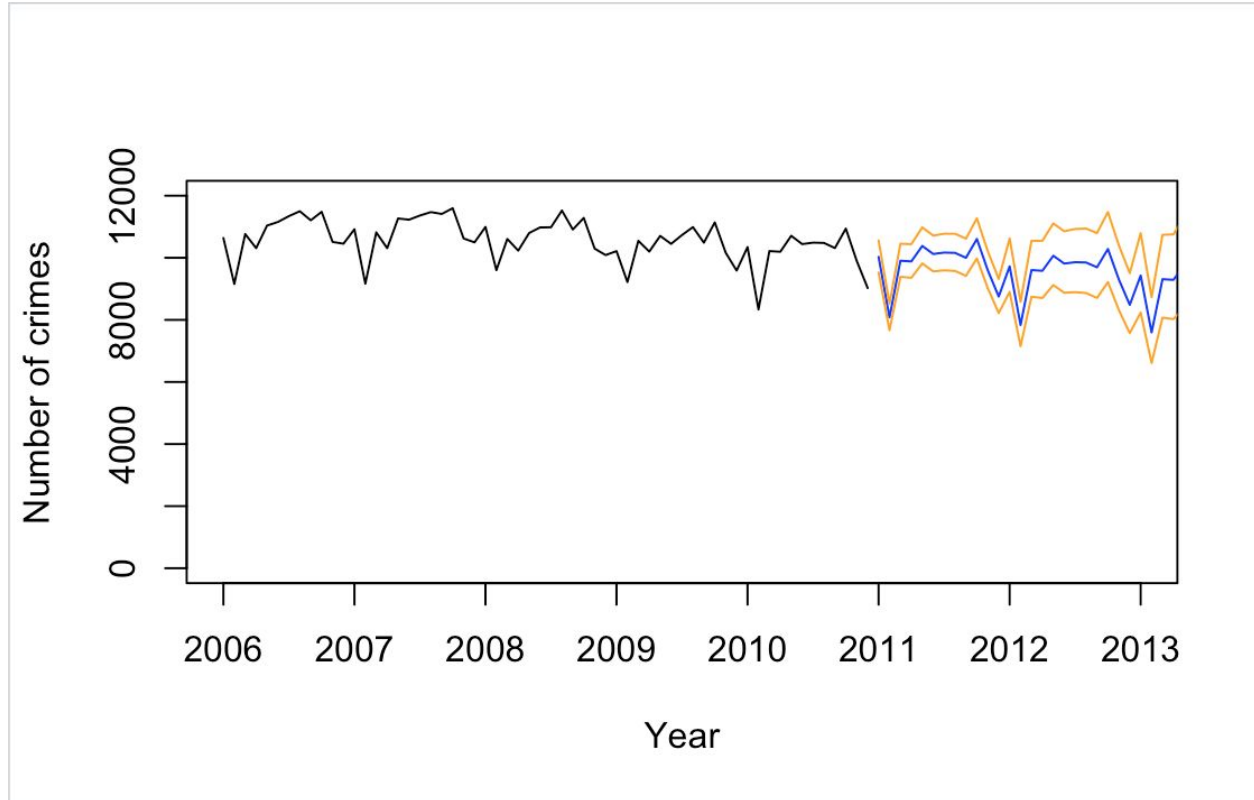- $q$ is the number of lagged forecast errors in the prediction equation.

# Crime time series analysis [original]

# Crime time series analysis [smoothing and difference]

# Crime time series analysis [forecasting]

# Practical considerations

**Missing values**
- ARIMA models and smoothing cannot be applied to time series with missing values
- Linear/logistic regression and neural network models do not require imputation

**Unequally spaced series**
- Many time series are naturally discrete (e.g. bus time arrivals, concerts, bid timings)
- Some forecasting methods might require interpolation

**Extreme values**
- Unusually large or small values in the series can affect forecasting
- Determine source of outliers (e.g. data entry errors, unusual events) and remove accordingly
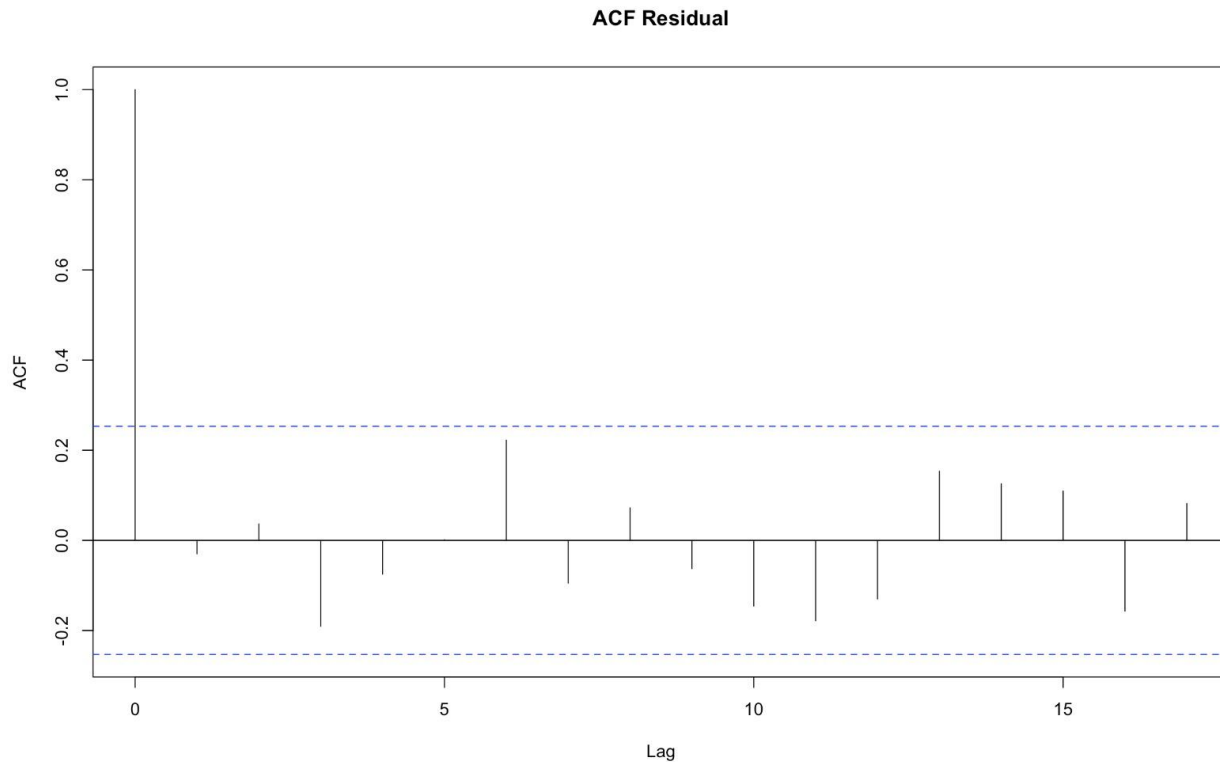
**Choice of time span**
- Long past of the series might deteriorate forecasting accuracy (changing context/environment)
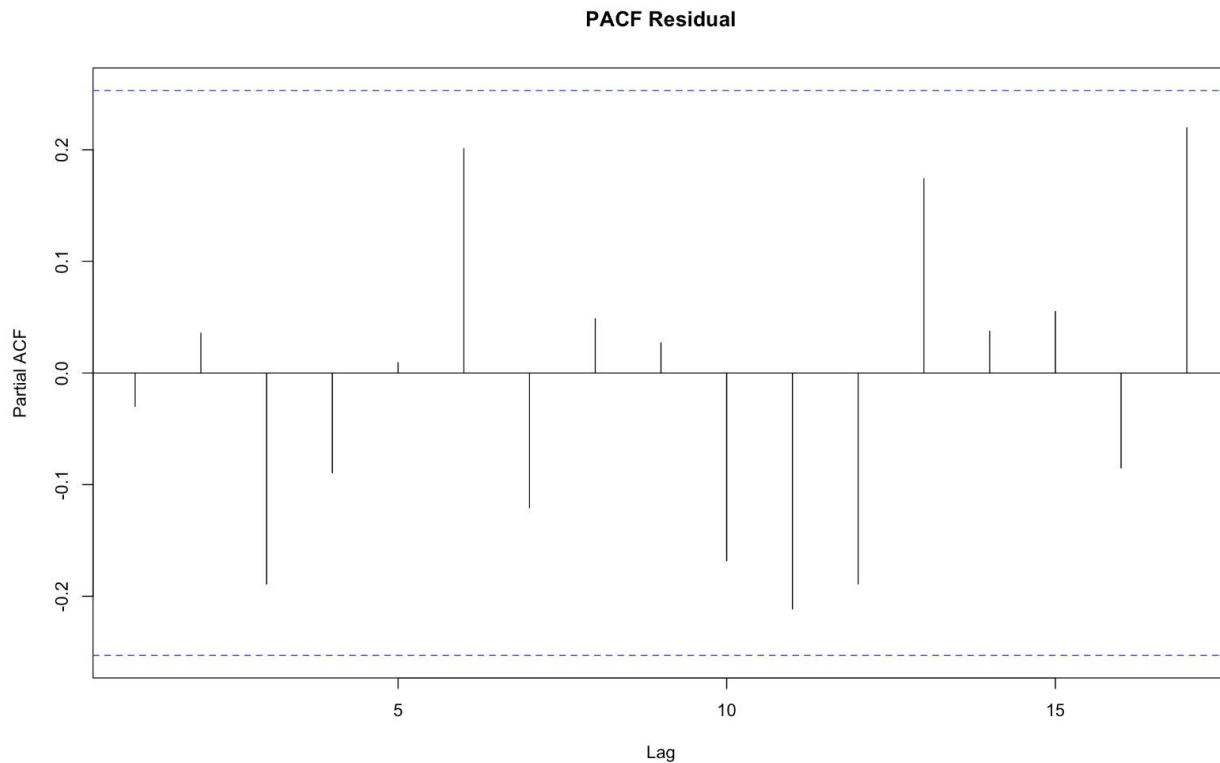
# Future work

- Verify fit of ARIMA model by examining residuals to ensure no more information is left for extraction (residuals should be random with no visible trends)

- Forecast heat maps for crime frequency at fine-grained spatial and temporal resolution

- Include other relevant data sources as exogenous variables to have correlated time series

- Classify crime types on time series, e.g. using long short-term memory (LSTM) with fuzzy logic

- Cluster time series of various different crime types based on individual crime type time series

**Questions?**

# ARIMA residuals [ACF diff log]



ACF Residual

# ARIMA residuals [PACF diff log]



PACF Residual

# Related work

J. R. Hipp, and C. E. Kubrin, "**Crime Report for Southern California**," UCI Irvine Laboratory for the Study of Space and Crime, 2015.

- Analyze the level of crime for all cities in the Southern California region with the population size of at least 4000
- Presents top/bottom cities for a specific crime type, and the increases/decreases of crime type
- Adjust the rates for the different socio-demographic characteristics of the city

H. Wang, D. Kifer, C. Graif, and Z. Li, "**Crime Rate Inference with Big Data**," Pennsylvania State University.

- Uses large-scale ("Big Data") Point-Of-Interest data as well as taxi flow data for Chicago
- Observed consistently improved performance in predicting crime rates for multiple years compared to using just demographic features

H. Kang, and H. Kang, "**Prediction of crime occurrence from multi-modal data using deep learning**," 2017.

- Uses feature-level data fusion method based on deep neural network (DNN)
  - Train DNN with spatial, temporal, environmental context, and joint feature representation layers
- Dataset consists of: crime statistics, demographic and meteorological data, and images of Chicago
- Improved accuracy in predicting crime