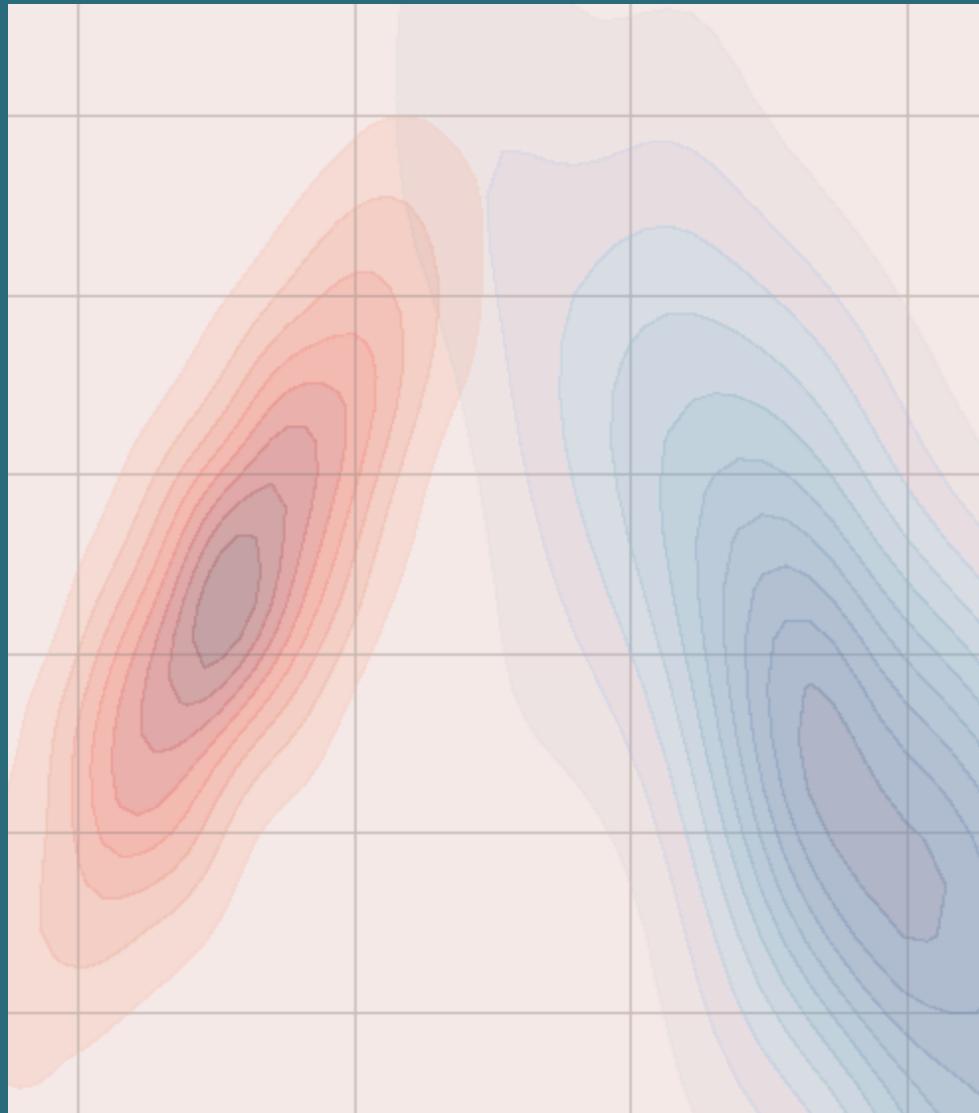


# Interrupted Time Series Experiments in Python

Drew Fustin  
Automaton Data  
[drewfustin@gmail.com](mailto:drewfustin@gmail.com)  
[@drewfustin](https://twitter.com/drewfustin)



# Me

**Automaton Data [contracting]**

drew@automatondata.com

PhD, Physics

Data Scientist



SPOT  
**HERO**  
**GRUBHUB**



@drewfustin | 2017.12.04

# Experiments are Hard



I HAVE NO  
IDEA WHAT  
I'M DOING

# Some are Straightforward



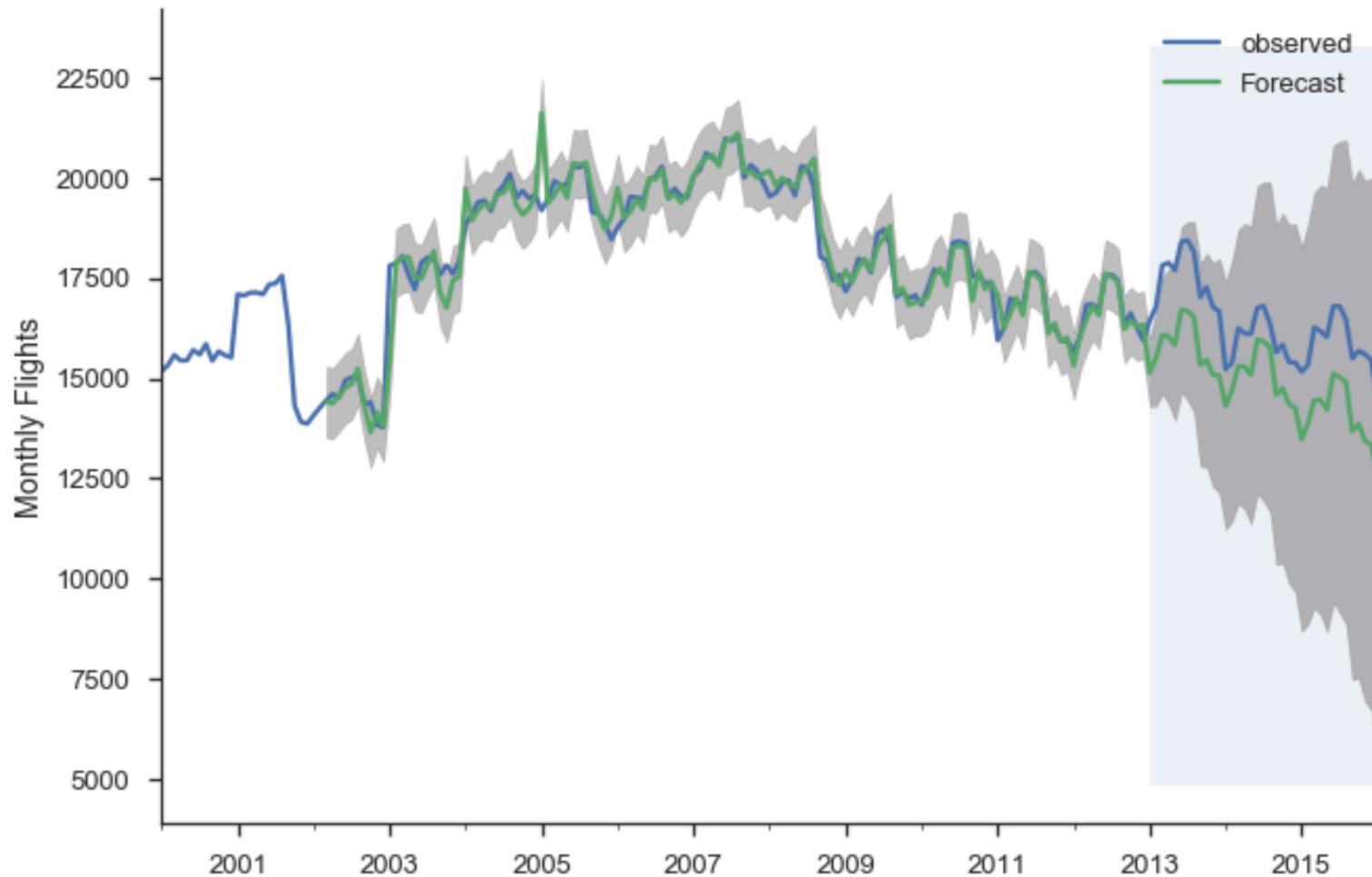
# Some are Straightforward

$$\Pr(p_B > p_A) = \sum_{i=0}^{\alpha_B - 1} \frac{B(\alpha_A + i, \beta_B + \beta_A)}{(\beta_B + i)B(1 + i, \beta_B)B(\alpha_A, \beta_A)}$$

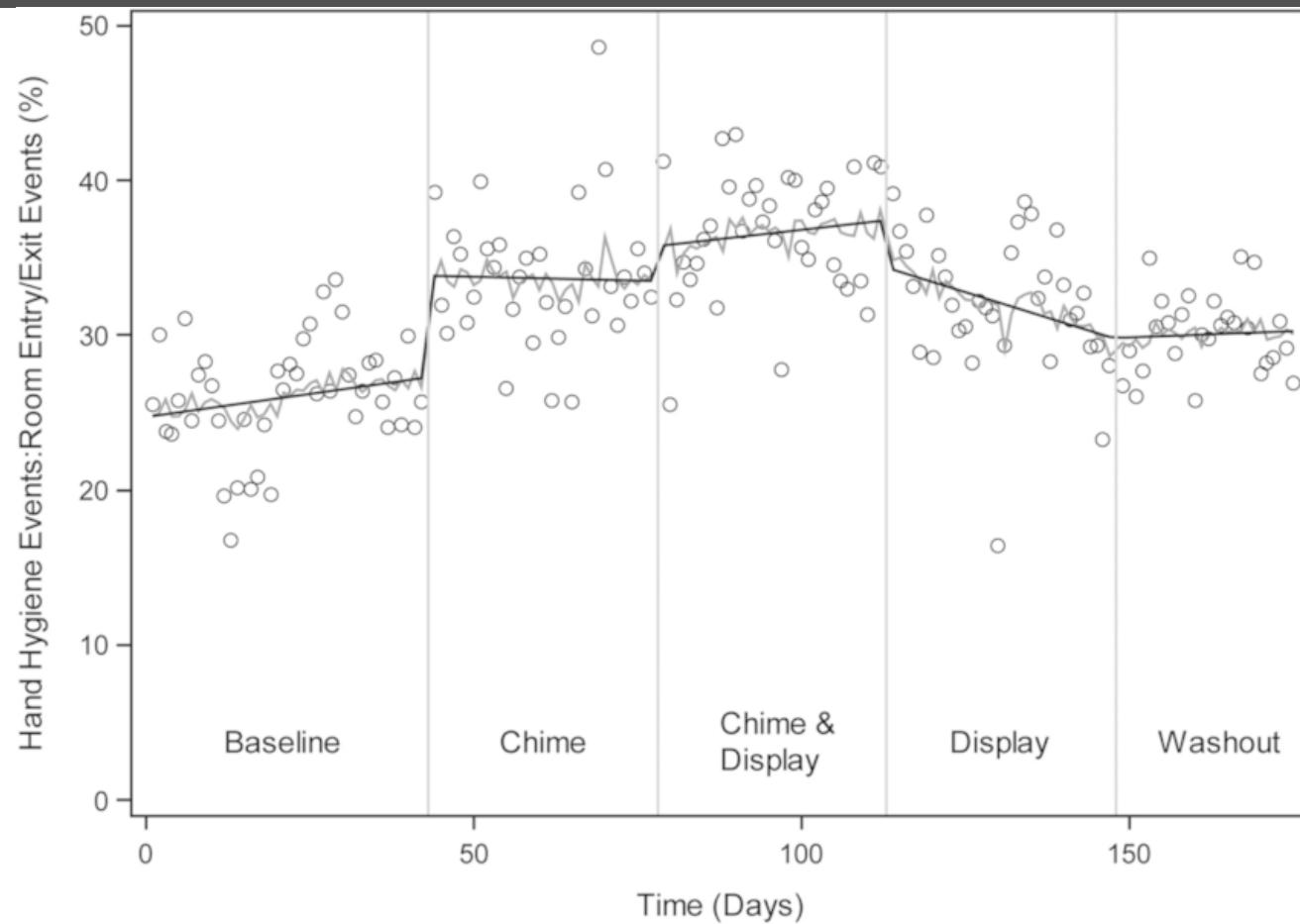
Where:

- $\alpha_A$  is one plus the number of successes for A
- $\beta_A$  is one plus the number of failures for A
- $\alpha_B$  is one plus the number of successes for B
- $\beta_B$  is one plus the number of failures for B
- $B$  is the [beta function](#)

# Some are Convolved

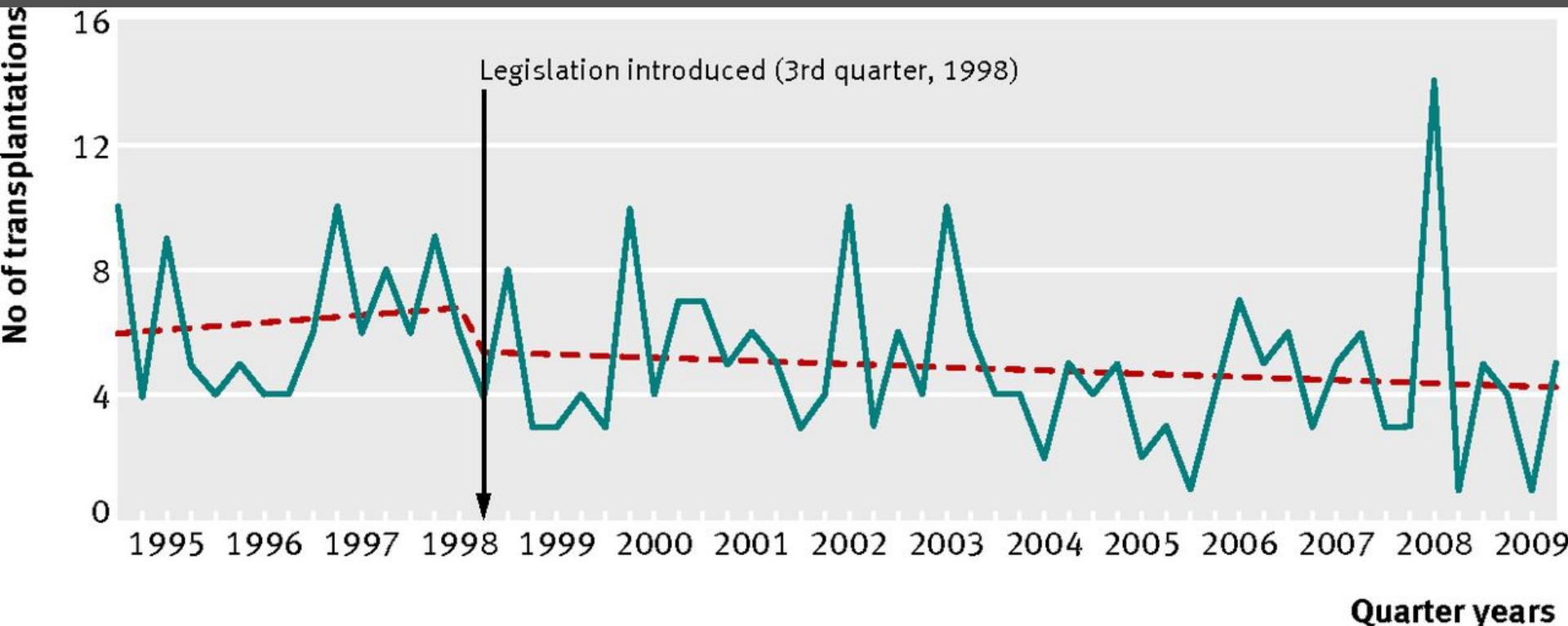


# Interrupted Time Series



Ellison, et al  
A Prospective Controlled Trial of an Electronic Hand  
Hygiene Reminder System

# Interrupted Time Series



Hawton, et al

Long term effect of reduced pack sizes of paracetamol on  
poisoning deaths and liver transplant activity in England and  
Wales: interrupted time series analyses

# Interrupted Time Series

So shaky, they're called  
**quasi-experiments**

# Interrupted Time Series

Related to

**regression discontinuity**

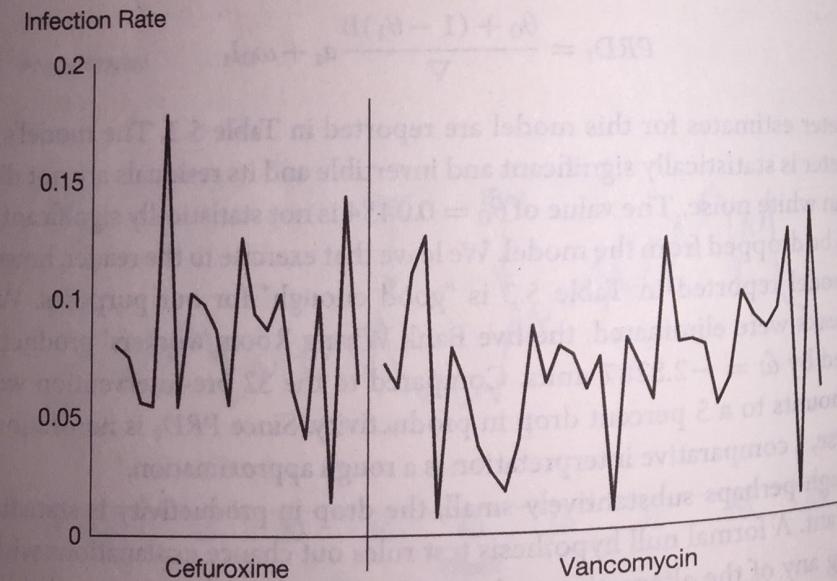
as time is the running variable,  
it is not randomly assigned  
*(required of RD design)*

# Interrupted Time Series

basic idea #1: fit ARIMA + effect size

for the next  $N_{post} = 35$  months, with prophylactic vancomycin. Modeling time series will be especially difficult. The inherent difficulty of identifying  $I_t$ ) from a short time series is aggravated in this instance by the unruly appearance of this time series and the possibility of a large impact. We start the alternative ARIMA modeling strategy with the simplest possible model:

$$SSI_t = \theta_0 + a_t + \omega_0 I_t$$

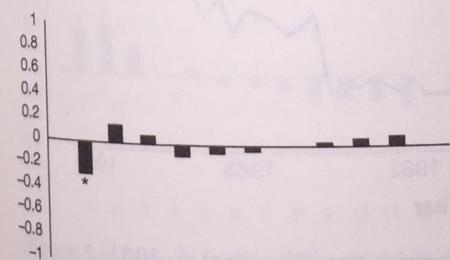


$$\begin{aligned} SSI_t &= 0.084 - 0.014I_t + \frac{a_t}{1 + .27B} \\ &= 0.084 - 0.014I_t + \sum_{k=0}^{\infty} 0.27^k a_t \\ &= 0.084 - 0.014I_t + a_t + 0.27a_{t-1} + 0.07a_{t-2} + 0.02a_{t-3} + \dots \end{aligned}$$

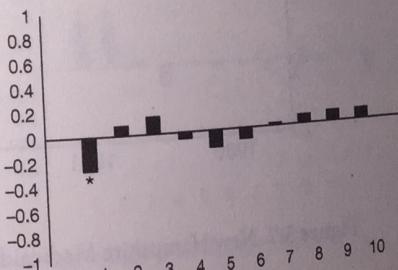
Table 5.4. SURGICAL SITE INFECTION RATES, ESTIMATION AND DIAGNOSIS

	Estimate	SE	t	Diagnosis	RSE
Model 1	$\hat{\theta}_0$	.0833	.0075	11.04	0.03373
	$\hat{\omega}_0$	-.0141	.0095	-1.49	$Q_{10} = 7.4$ $Z_{KS} = 0.805$
Model 2	$\hat{\theta}_0$	.0836	.0059	14.11	0.03277
	$\hat{\phi}_1$	-.2710	.1312	-2.07	$Q_{10} = 3.9$
	$\hat{\omega}_0$	-.0144	.0074	-1.95	$Z_{KS} = 0.805$

Sample ACF for  $\hat{a}_t$

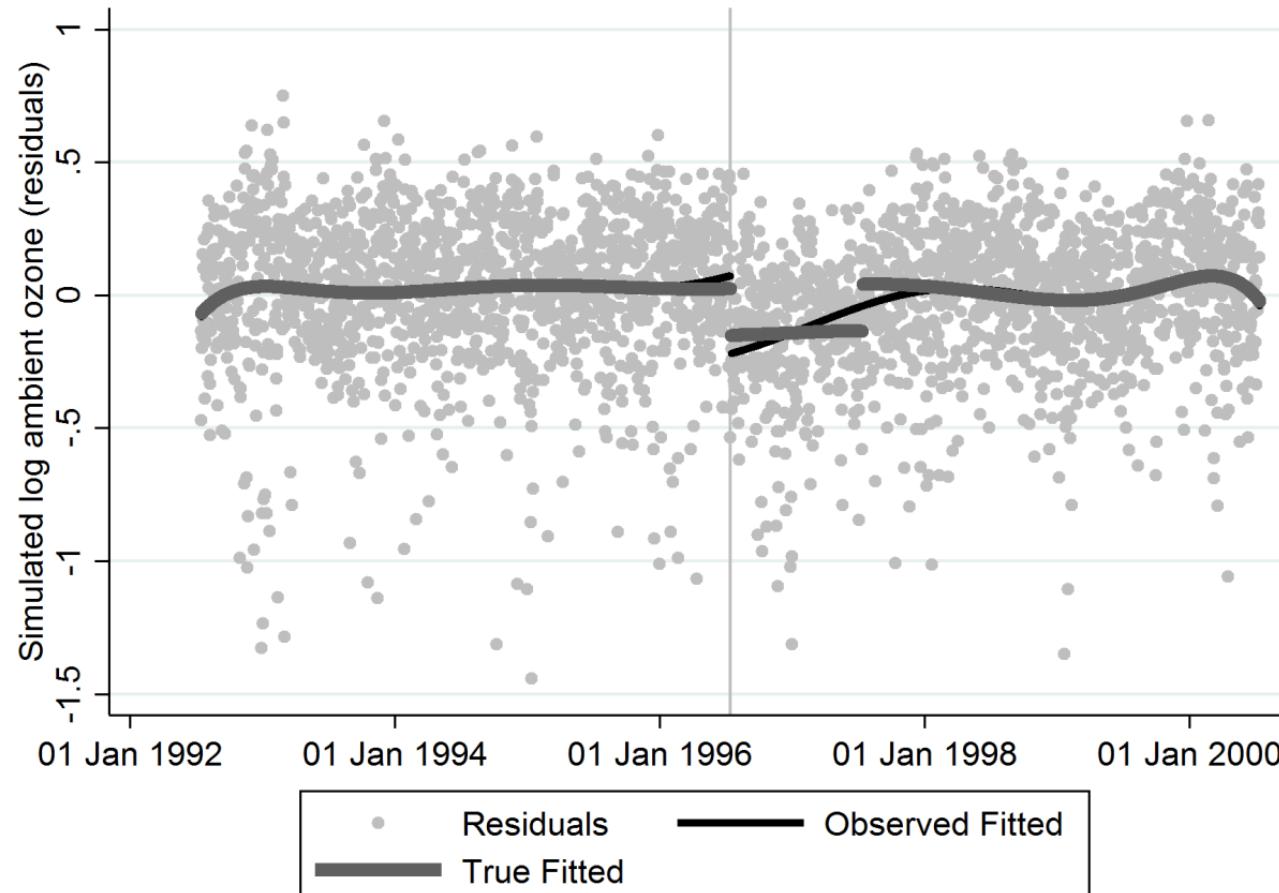


Sample PACF for  $\hat{a}_t$



# Interrupted Time Series

basic idea #2: fit curves at discontinuity



# Let's Simulate Some Data

## Homogeneous Poisson Point Process    $\lambda$ is fixed in time

- $N(0) = 0$ ;
- has independent increments; and
- the number of events (or points) in any interval of length  $t$  is a Poisson random variable with parameter (or mean)  $\lambda t$ .

The last property implies:

$$E[N(t)] = \lambda t.$$

In other words, the probability of the random variable  $N(t)$  being equal to  $n$  is given by:

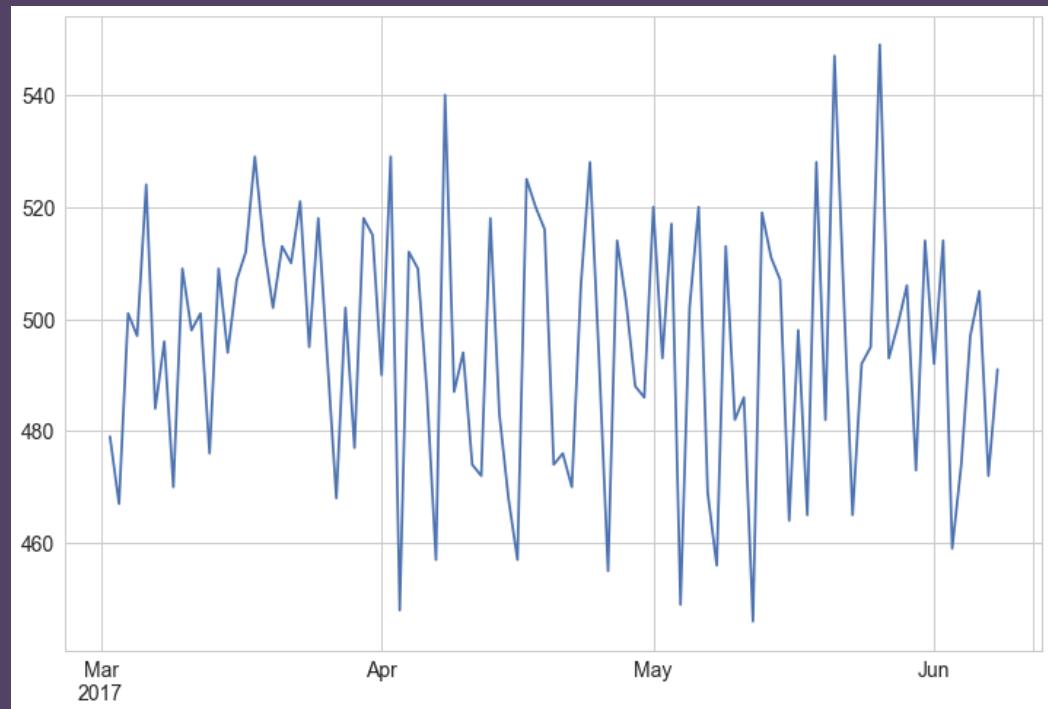
$$P\{N(t) = n\} = \frac{(\lambda t)^n}{n!} e^{-\lambda t}.$$

The Poisson counting process can also be defined by stating that the time differences between events of the counting process are exponential variables with mean  $1/\lambda$ .<sup>[52]</sup> The time differences between the events or arrivals are known as **interarrival** <sup>[53]</sup> or **interoccurrence** times.<sup>[52]</sup>

# Let's Simulate Some Data

## Homogeneous Poisson Point Process    $\lambda$ is fixed in time

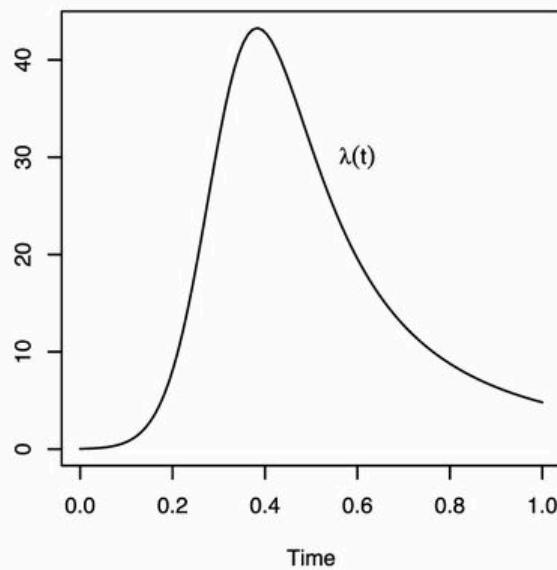
cumulatively  
add events  
distributed  
exponentially  
with rate  $1/\lambda$



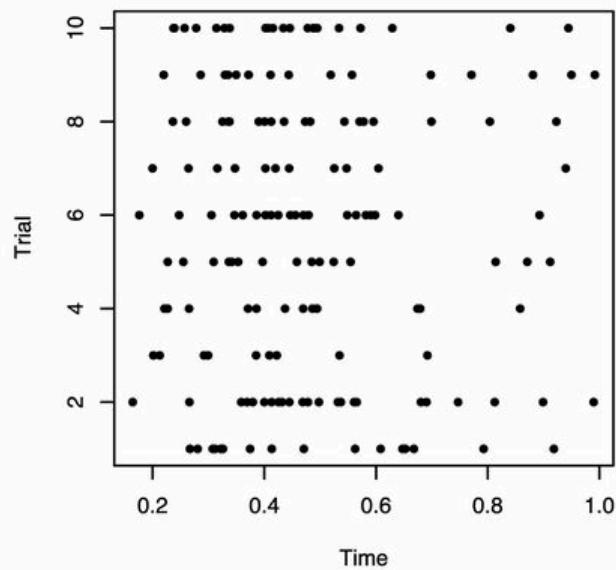
# Let's Simulate Some Data

## Nonhomogeneous Poisson Point Process $\lambda = \lambda(t)$

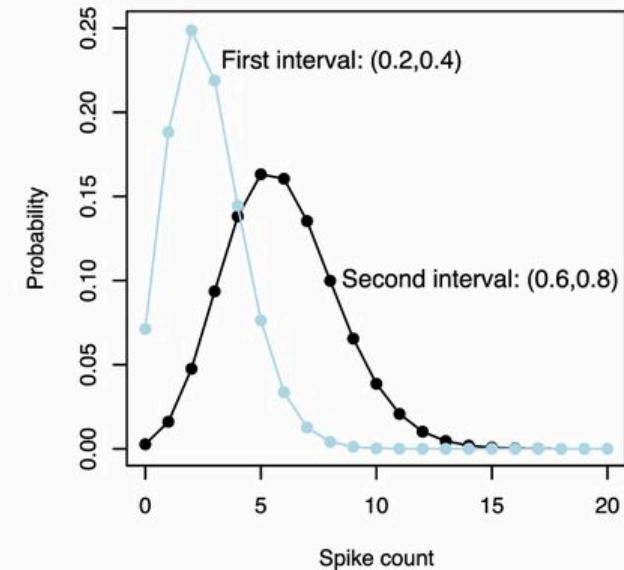
a.



b.

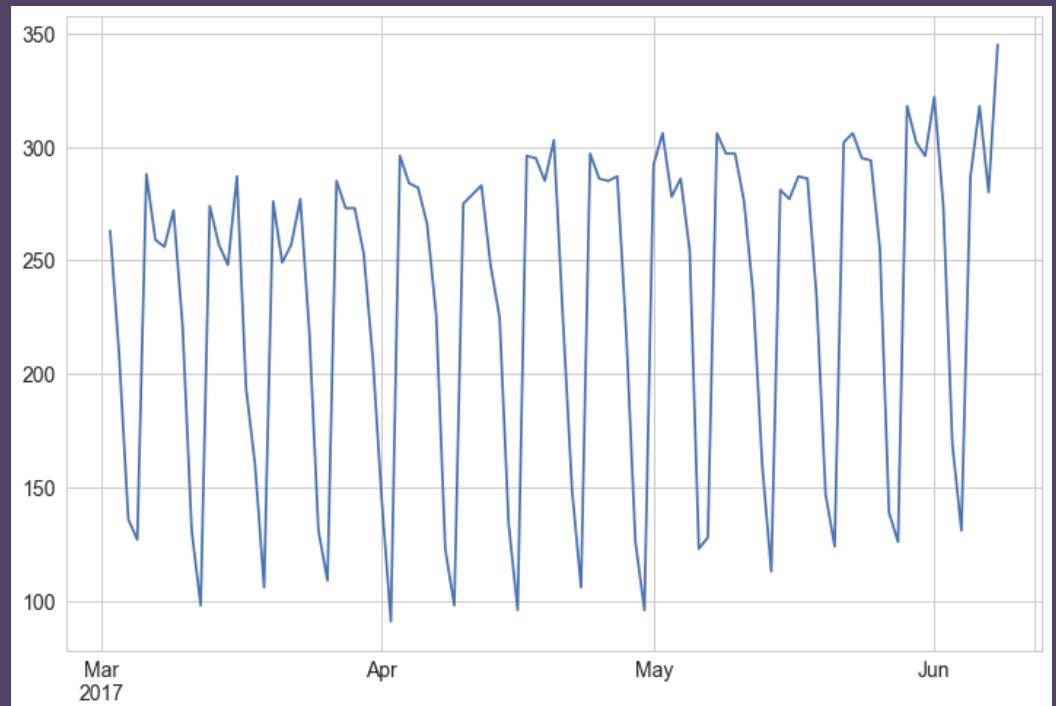
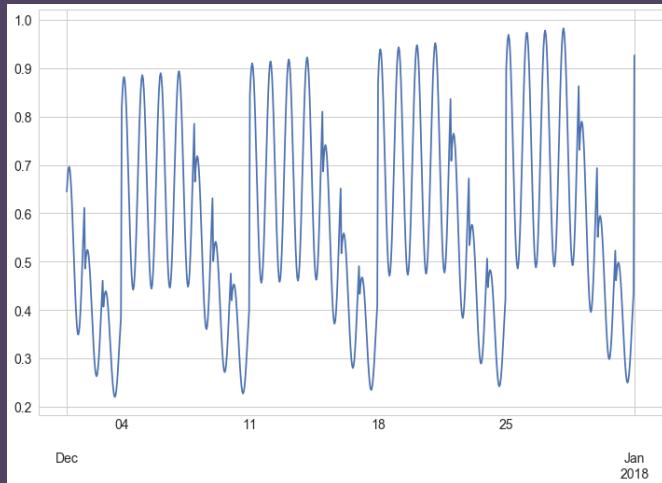


c.



# Let's Simulate Some Data

## Nonhomogeneous Poisson Point Process $\lambda = \lambda(t)$



# ARIMA Forecasting

Autoregressive (AR)

Integrated (I)

Moving Average (MA)

The above can be generalized as follows.

$$\left(1 - \sum_{i=1}^p \phi_i L^i\right) (1 - L)^d X_t = \delta + \left(1 + \sum_{i=1}^q \theta_i L^i\right) \varepsilon_t.$$

This defines an ARIMA( $p,d,q$ ) process with **drift**  $\delta/(1 - \sum \phi_i)$ .

# ARIMA Forecasting



# ARIMA Forecasting

Autoregressive (AR)

$$\phi_p(B) = (1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p)$$

Integrated (I)

$$\theta_q(B) = (1 - \theta_1 B - \theta_2 B^2 - \dots - \theta_q B^q)$$

Moving Average (MA)

$$\nabla^d(B) = (1 - B)^d$$

e.g. ARIMA (1, 1, 0):

$$(1 - \phi_1 B)(1 - B) = 1 - (1 + \phi_1)B + \phi_1 B^2$$

# ARIMA Forecasting

$$\text{ARIMA } \underbrace{(p, d, q)}_{\begin{pmatrix} \text{Non-seasonal part} \\ \text{of the model} \end{pmatrix}} \quad \underbrace{(P, D, Q)_m}_{\begin{pmatrix} \text{Seasonal part} \\ \text{of the model} \end{pmatrix}}$$

period. For example, an ARIMA(1,1,1)(1,1,1)<sub>4</sub> model (without a constant) is for quarterly data ( $m = 4$ ) and can be written as

$$(1 - \phi_1 B) (1 - \Phi_1 B^4) (1 - B) (1 - B^4) y_t = (1 + \theta_1 B) (1 + \Theta_1 B^4) e_t.$$

Diagram illustrating the decomposition of the ARIMA(1,1,1)(1,1,1)<sub>4</sub> model:

- (1 -  $\phi_1 B$ )  $\xrightarrow{\text{(Non-seasonal)}} \text{AR}(1)$
- ( $1 - \Phi_1 B^4$ )  $\xrightarrow{\text{(Seasonal)}} \text{AR}(1)$
- (1 -  $B$ )  $\xrightarrow{\text{(Non-seasonal)}} \text{difference}$
- (1 -  $B^4$ )  $\xrightarrow{\text{(Seasonal)}} \text{difference}$
- (1 +  $\theta_1 B$ )  $\xrightarrow{\text{(Non-seasonal)}} \text{MA}(1)$
- (1 +  $\Theta_1 B^4$ )  $\xrightarrow{\text{(Seasonal)}} \text{MA}(1)$

# ARIMA Forecasting



# ARIMA Forecasting

## Other Topics

- Stationarity: Adjusted Dickey-Fuller Test
- Autocorrelation Functions: ACF/PACF
  - reading *those* tea leaves
- “Memory” differences between AR, I, MA

# PyMC3 for ITS Analysis

Generate NHPP

Fit ARIMA to get (p, d, q) order

Add response  $\sim \mathcal{N}(\mu_\beta, \sigma_\beta)$  at t\_exp

Use ARIMA ops to iteratively solve y,  $\varepsilon$

# PyMC3 for ITS Analysis

Generate NHPP

Fit ARIMA to get (p, d, q) order

Add response  $\sim \mathcal{N}(\mu_\beta, \sigma_\beta)$  at t\_exp

Use ARIMA ops to iteratively solve y,  $\varepsilon$

Potential next step:

Use PyMC3 to check regression  
discontinuity in across t\_exp bound