

3 Epidemic Data and Time Series Tools



Mathematics
and Statistics

$$\int_M d\omega = \int_{\partial M} \omega$$

Mathematics 4MB3/6MB3 Mathematical Biology

Instructor: David Earn

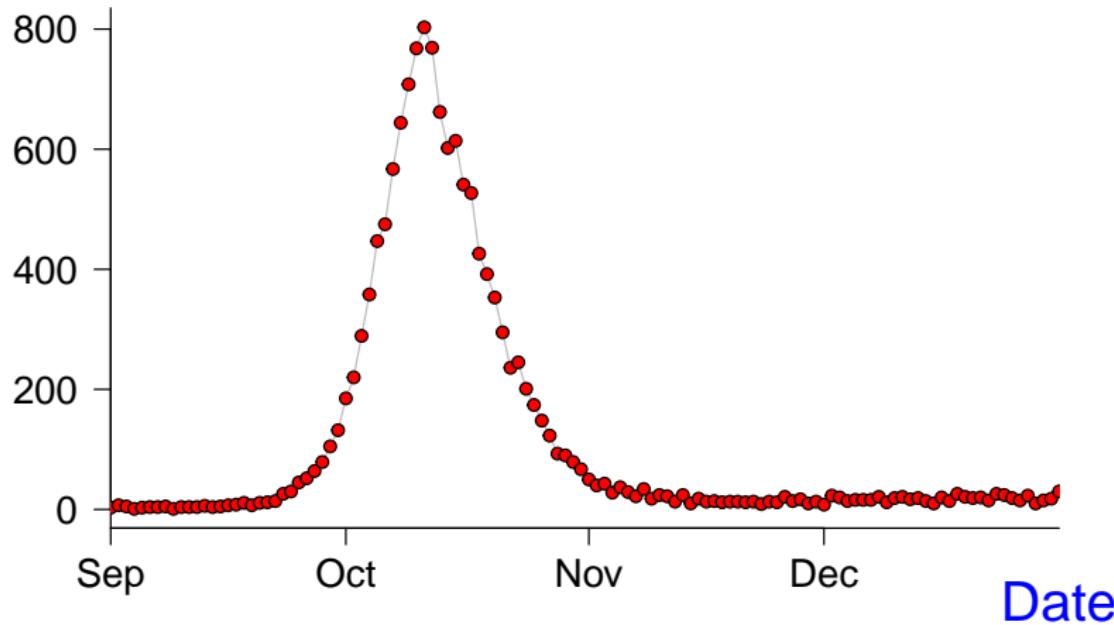
Lecture 3
Epidemic Data and Time Series Tools
Tuesday 17 September 2024

Announcements

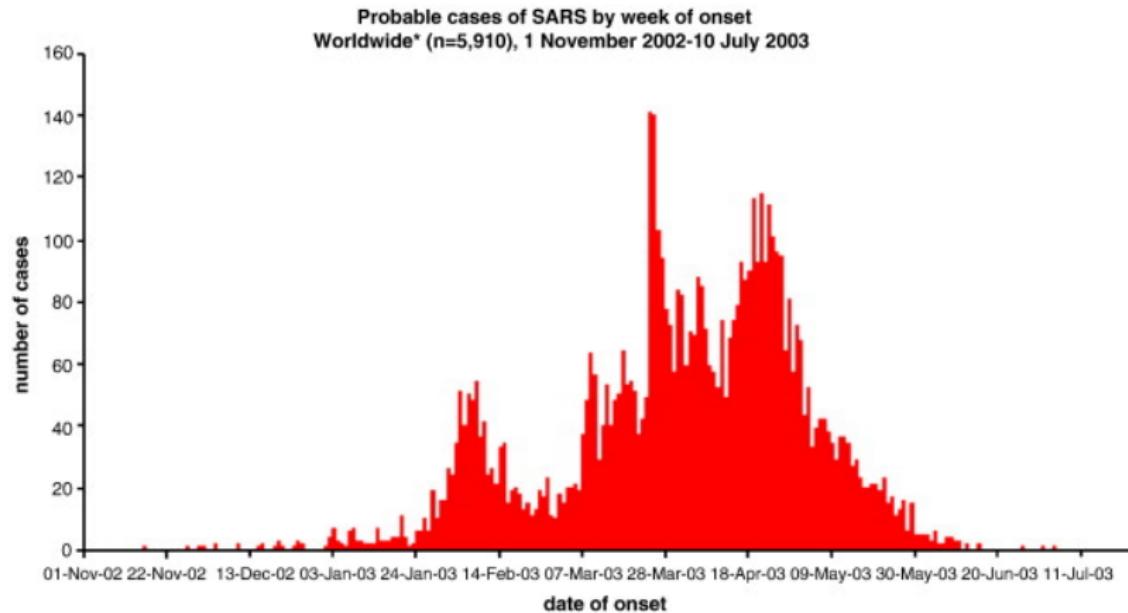
- Next week's lecture will be recorded in advance and posted on the Echo 360 page for this course.
 - Live Q&A, either in last hour of scheduled class or at a mutually convenient time.
- Assignments:
 - Assignment 1 due 23 Sep 2024 (next Monday)
 - Assignment 2 due 7 Oct 2024
(good to work on before class on 1 Oct 2024)
- Class on 1 Oct 2024 will be given by Mikael Jagan
(install epigrowthfit before that class)
- Lecture on 8 Oct 2024 will pre-recorded and posted

P&I Mortality, Philadelphia, 1918

P&I Deaths

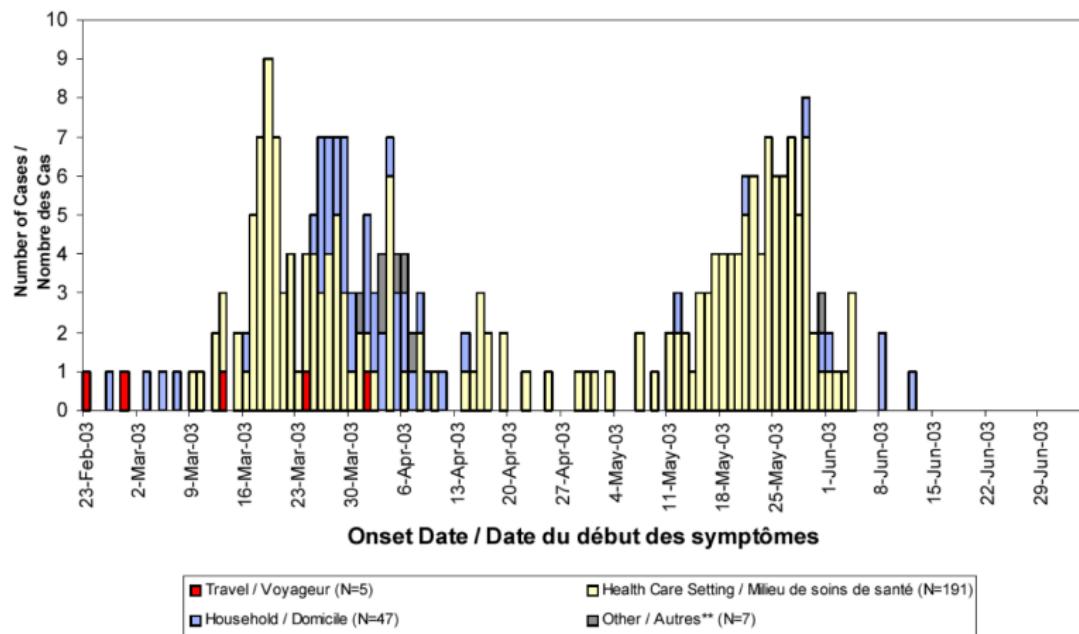


SARS in 2003 (Worldwide)



*This graph does not include 2,527 probable cases of SARS (2,521 from Beijing, China), for whom no dates of onset are currently available.

SARS in 2003 (Toronto)



$N = 249$ (of 250 reported)

Some SARS Facts

- High case fatality
 - 1918 flu < 3%
 - SARS > 10%
- Long hospital stays
 - Mean time from admission to discharge or death:
~ 25 days in Hong Kong
- 8098 probable cases, 774 deaths
- How bad would it have been if it had not been controlled?

COVID (ancestral) hospitalization and survival in Ontario

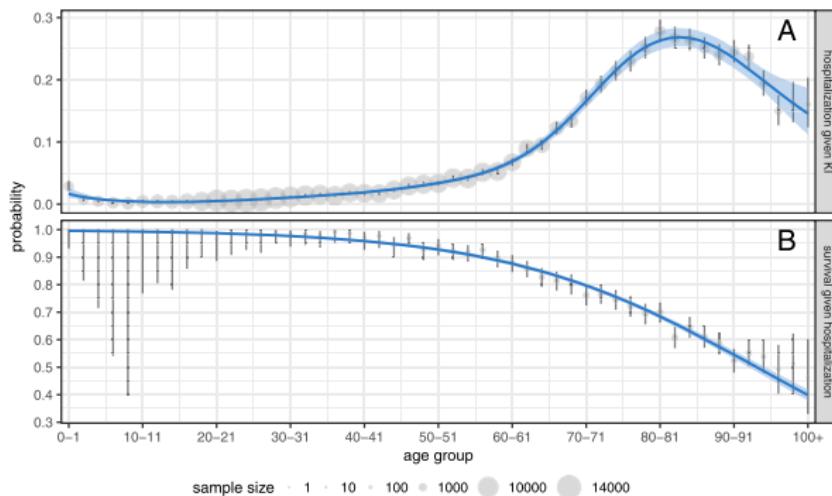
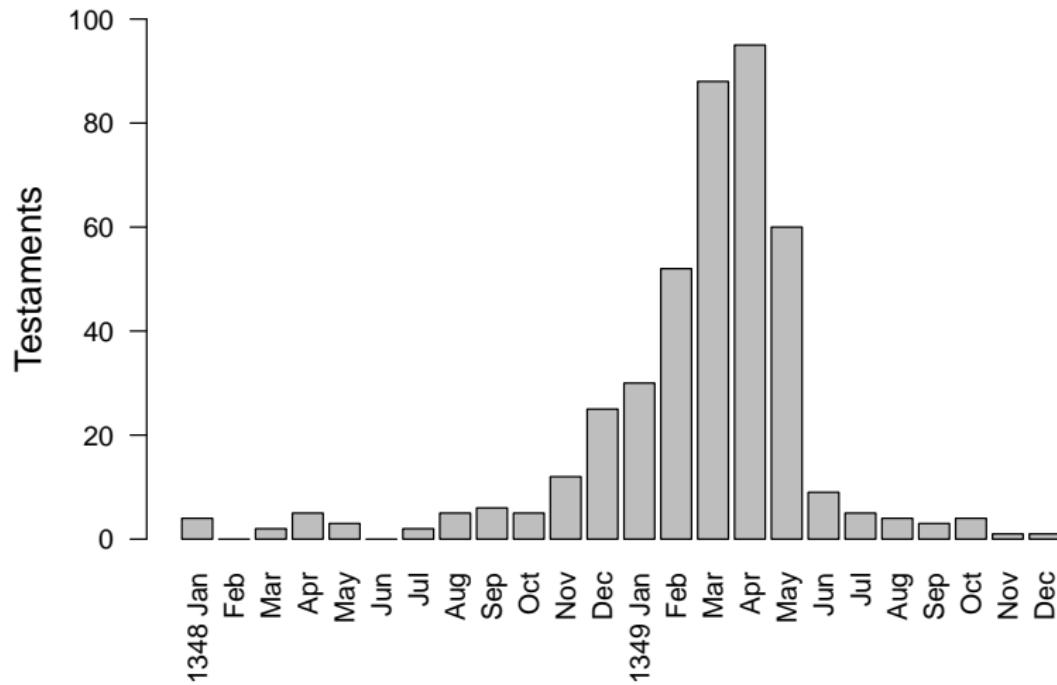


Fig. 4 Age-dependent COVID-19 hospitalization probability for known SARS-CoV-2 infection (panel **a**) and survival probability for hospitalized patients (panel **b**) in Ontario. We give age-by-age estimates of each probability (points; 95% exact binomial confidence intervals given by vertical lines), where point area is proportional to age-specific sample size. We additionally provide more precise estimates of these probabilities under stricter assumptions, modelling the hospitalization probability using a generalized additive model and the survival probability using a generalized linear model (curves; 95% confidence bands given by shaded regions). See “Methods” section for details

Papst et al (2021), BMC Public Health

The Black Death in London, England, 1348–1349



London Bill of Mortality, 26 Sept to 3 Oct 1665

London Bill of Mortality, 26 Sept to 3 Oct 1665

Frighted	
Gowt	1
Grief	1
Griping in the Guts	3
Jaundies	35
Imposthume	2
Infants	8
Kingsevil	9
Mcagrome	2
Plague	5533
Purples	2
Ricketts	

Mortality Bills are typically handwritten

LONDON		From the 4 th of July to the 11 th of the same 1665.					
Buried	Plag.	Buried	Plag.	Buried	Plag.	Buried	Plag.
St Alban Woodstreet	2	St Clement Eastchapel	1	St Margaret Newfisht		St Michael Crookedst.	4
Alhallows Bark-	2	St Dionys Backchurch	1	St Margaret Patrons		St Michael Queenhith	3
Alhallows Breadstreet		St Dunstan East	2	St Mary Abchurch		St Michael Quern	
Alhallows Great	1	St Edmund Lombardst.	2	St Mary Aldermanbury		St Michael Royal	
Alhallows Honouranc-		St Ethelborrough		St Mary Aldemary		St Michael Woodstreet	
Alhal'ows Less	1	St Faiths	1	St Mary le Bow		St Mildred Breadstreet	
Alhallows Lombardstr		St Gabriel Fenchurch		St Mary Botham		St Mildred Poultrey	
Alhallows Staining	1	St George Borophlane		St Mary Colechurch		St Nicholas Acons	
Alhallows the Wall	1	St Gregoryes by St. Paul	3	St Mary Hill		St Nicholas Colcabby	
St Alphage		St Hellen	2	St Mary Mag. Milkstr.		St Nicholas Olaves	
St Andrew Hubbard		St James Dukes place	1	St Mary Mag. Oldfisht		St Olave Hartfreet	
St Andrew Underhaft	3	St James Garlickhithe	1	St Mary Mounthaw		St Olave Jewry	
St Andrew Wardrobe		St John Bapstif		St Mary Summerle	2	St Olave Silverstreet	4
St Anne Aldersgate		St John Evangelist		St Mary Staining		St Pancras Soperlane	1
St Anne Blackfriars	7	St John Zichary		St Mary Woolchurch		St Peter Cheap	
St Ancholins Parish		St Katharine Coleman		St Mary Woolnoth		St Peter Cornhil	
St Austin's Parish		St Katharine Creechar.		St Martins Iremongerl.		St Peter Paulwharf	
St Barthol. Exchange		St Lawrence Jewry		St Martins Ludgate	2	St Peter Poor	1
St Bennet Fynck		St Lawrence Pountney		St Martins Orgars		St Steven Colemaistr.	2
St Benne Gracechurch	7	St Leonard Eastchapel		St Martins Outwich	1	St Steven Walbrook	1
St Benner Paulswarf		St Leonard Fosterlane		St Martins Vintrey	1	St Switthin	2
St Benner Sherehog		St Magnus Parish	1	St Matthew Frydaystr.		St Thomas Apostle	1
St Borolph Billingsgate		St Margaret Lothbury		St Michael Bassiflsh	4	Trinity Parish	1
Christ Church		St Margaret Moles		St Michael Cornhil		St Vedast alias Fosters	
St Christopher's	6	Christened in the 9 th Parishes within the walls		Buried	86	Plague	28
St Andrew Holborn	66	St Borolph Aldersgate	11	St George Southwark	13	St Sepulchres Parish	117
St Bartholomew Great	40	St Borolph Aldgate	24	St Giles Cripplegate	103	St Thomas Southwark	81
St Bartholomew Less	4	St Borolph Bishopgate	37	St Olave Southwark	20	Trinity Minories	7
St Bridge	24	St Dunstan West	19	St Saviour Southwark	21	At the Pesthouse	6
Bridewel Preinct	14	Christened in the 15 th Parishes without the walls		Buried	473	Plague	243
Christ Church		St Kath. near the Tower	7	St Mary Fislington	3	St Paul Shadwei	
St John at Hackney	1	Lambeth Parish	4	St Mary Newington	7	Rotherhithe Parish	7
St Giles in the Fields	268	St Leonard Shoreditch	21	St Mary Whitechappel	16	Stepney Parish	47
St James Clerkenwel	213	St Magdalens Bermond.	13	Christened in the 12 th Out-Parishes in Middlesex and Surrey		Buried	455
	53		14			Plague	286

But handwriting is usually very clear

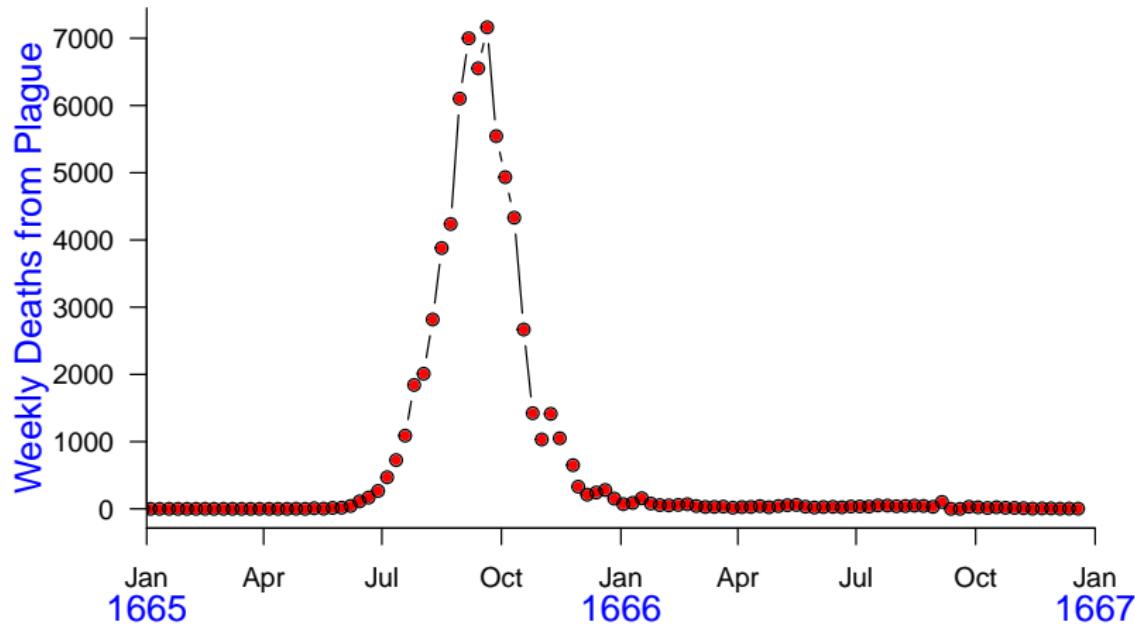
A historical ledger page from London, dated 29th [unclear]. The page is divided into columns for 'Buried.' and 'Plag.'. The data is organized by parish:

	Buried.	Plag.
St A lban Woodstreet	2	1
Alhallows Bark.-	2	
Alhallows Breadstreet	1	
Alhallows Great —	1	

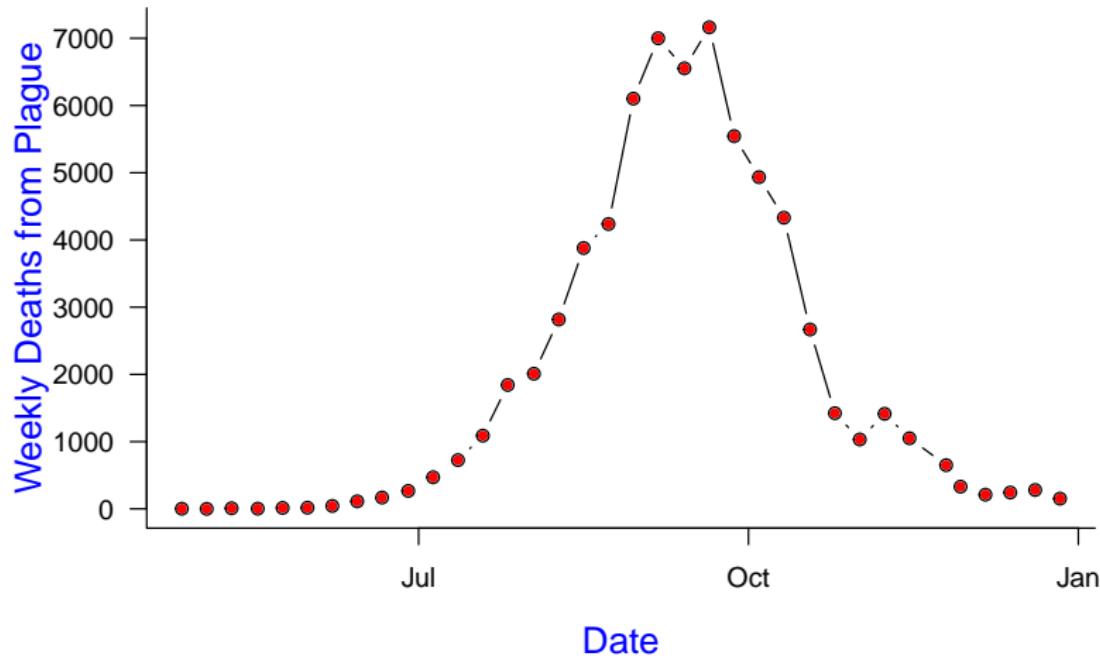
But handwriting is usually very clear

St Christopher's ———			
Christened in 97 the Parishes :			
St Andrew Holborn —	66	40	Se
St Bartholomew Great	+	+	Se
St Bartholomew Less —	+	+	Se
St Bridget ——— —	24	17	Se
Bridewell Precept —	1	1	Se
Christened in the 15 Parishes :			

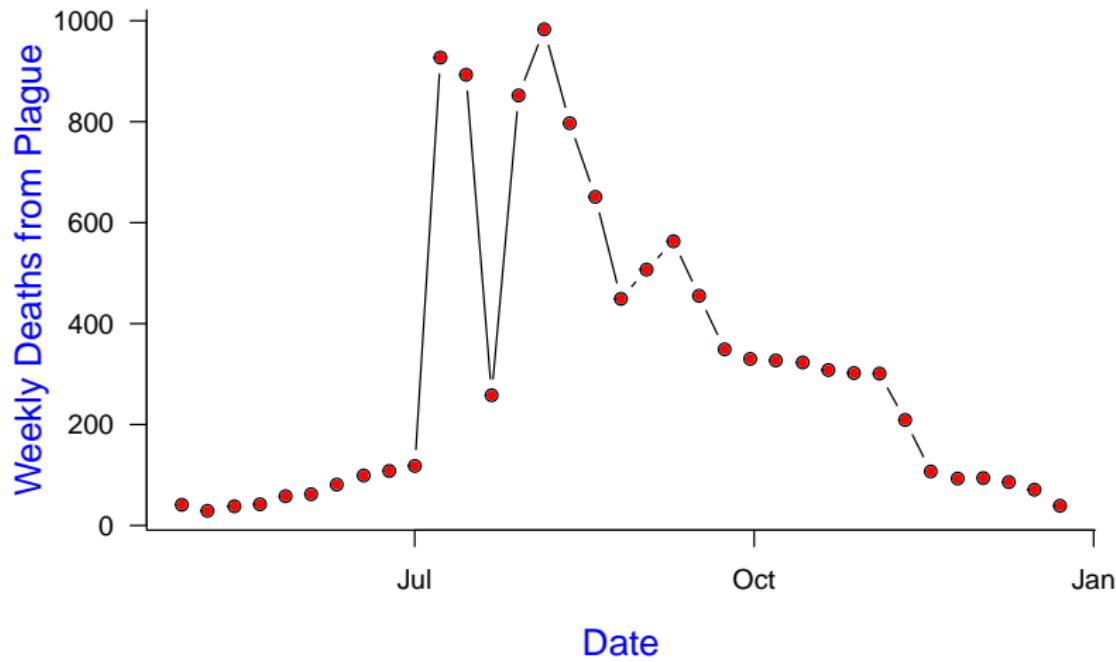
The Great Plague of London, 1665



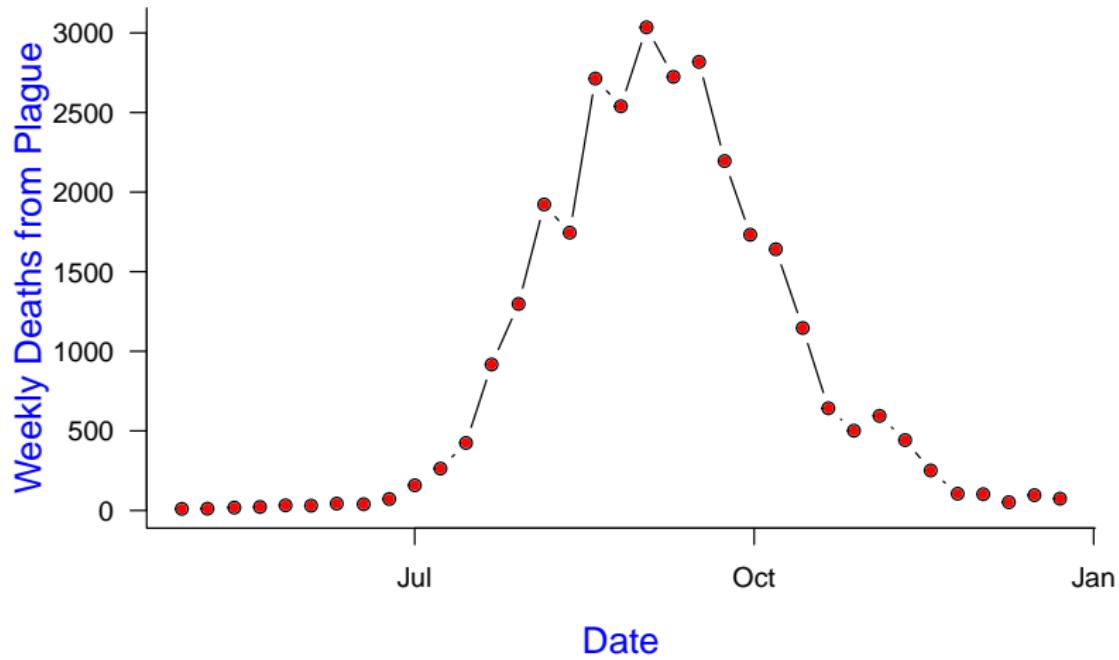
The Great Plague of London, 1665



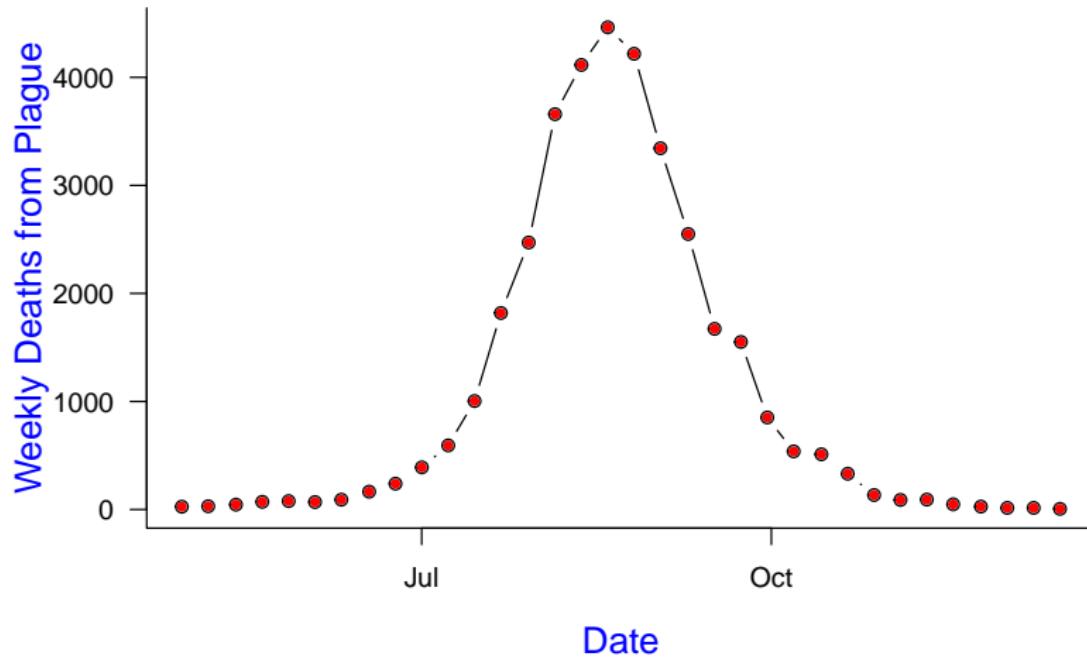
London Plague of 1593



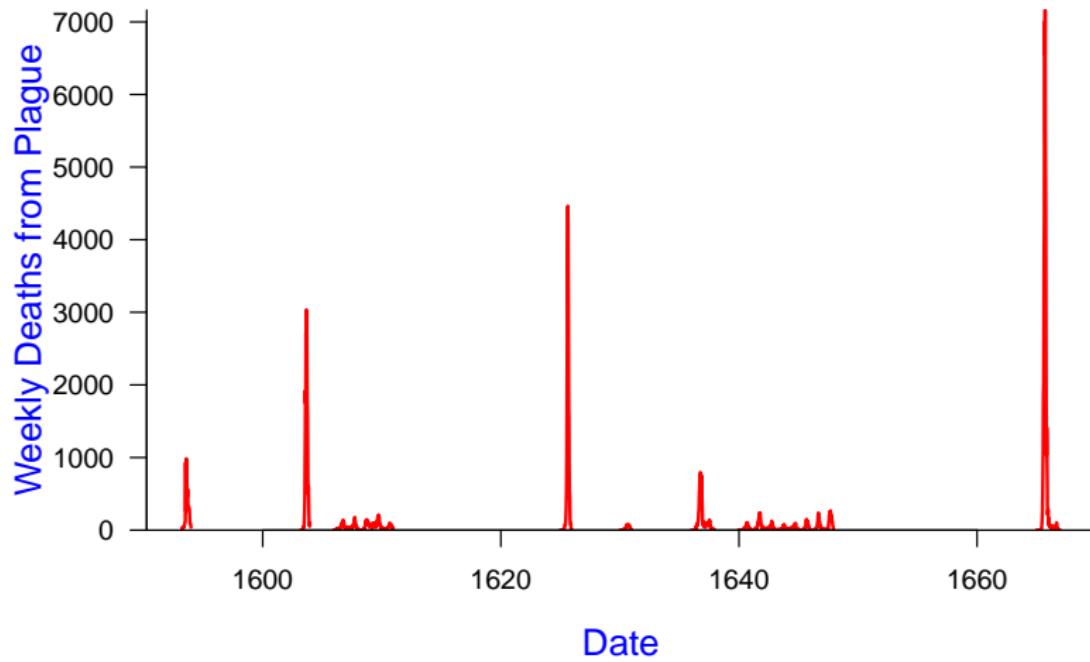
London Plague of 1603



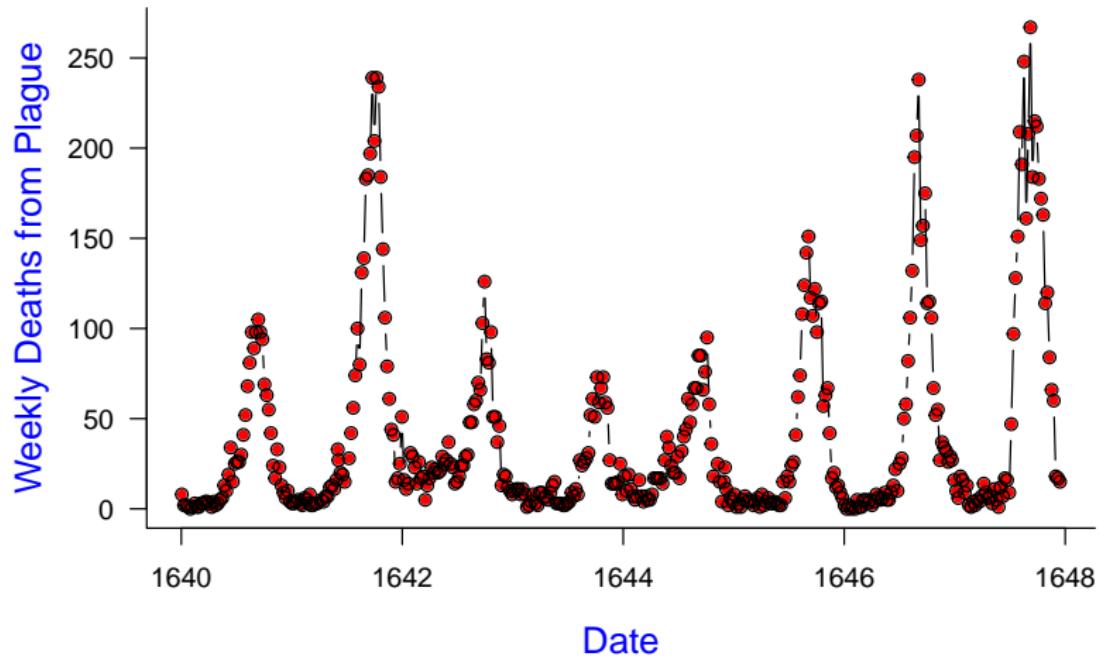
London Plague of 1625



Weekly Deaths from Plague in London, 1592–1666



Weekly Plague in London, 1640–1648



Some Plague Facts

- Plague epidemics recorded from Roman times to early 1900s.
- $\gtrsim 1/3$ Europe's population died in "Black Death" of 1348
 - ~ 300 years for the population to reach the same level.
- Recently (2011) established (at McMaster!) that the pathogen that caused The Black Death was *Yersinia pestis*

[Bos et al. 2011, *Nature* 478, 506–510]

- More recently (2014) established (again at McMaster!) that the pathogen that caused The Plague of Justinian (541–543 AD) was *Yersinia pestis*

[Wagner et al. 2014, *Lancet Infectious Diseases* 14, 319–326]

- *Y. pestis* still a concern?
Yes: Rodent reservoir, antibiotic-resistant strains, bioterrorism
- **Spatial data** for any plagues? Yes, for London in 1665...

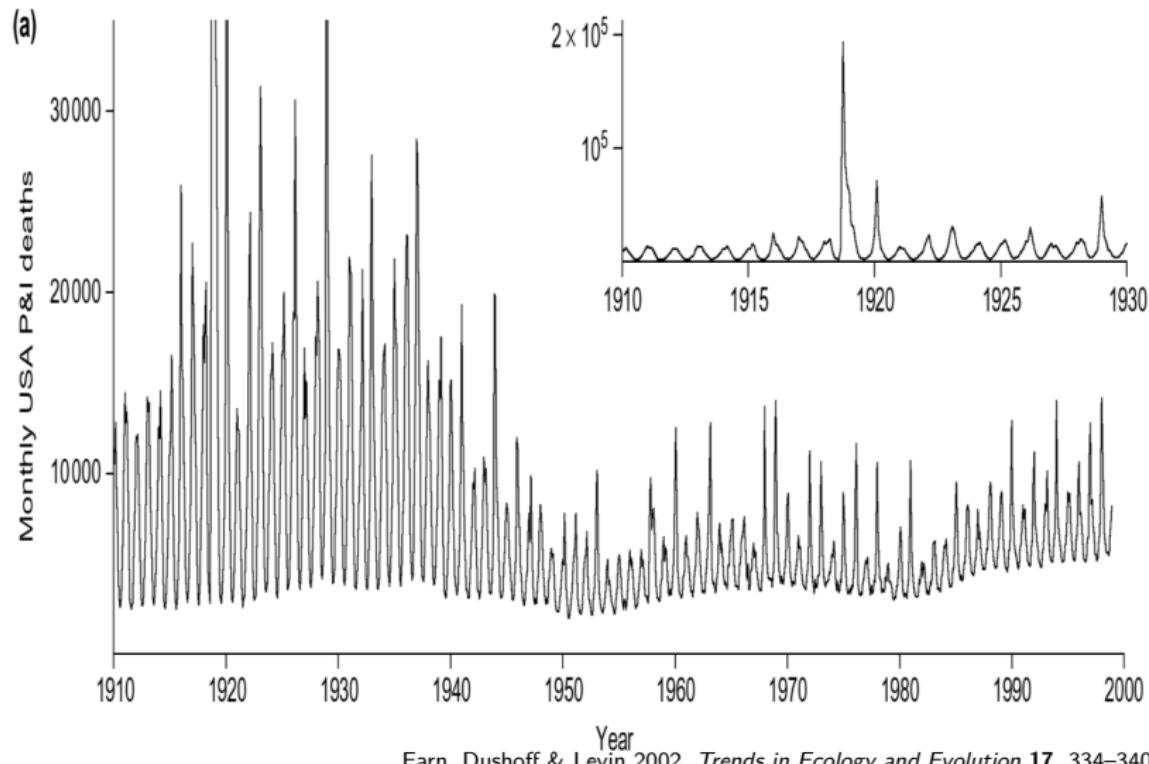
Visualization of spatial structure of Great Plague

- GIS encoding of parish boundaries
- Overlay parish boundaries on more modern map for reference
- Colour parishes as they become infected
- Is there evidence for spatial spread or was the spatial pattern random?
- DE low-tech animation...
- CBC high-tech animation...
 - *The Nature of Things*, 21 August 2014.
[http://www.cbc.ca/natureofthings/episodes/
secrets-in-the-bones-the-hunt-for-the-black-death-killer](http://www.cbc.ca/natureofthings/episodes/secrets-in-the-bones-the-hunt-for-the-black-death-killer)

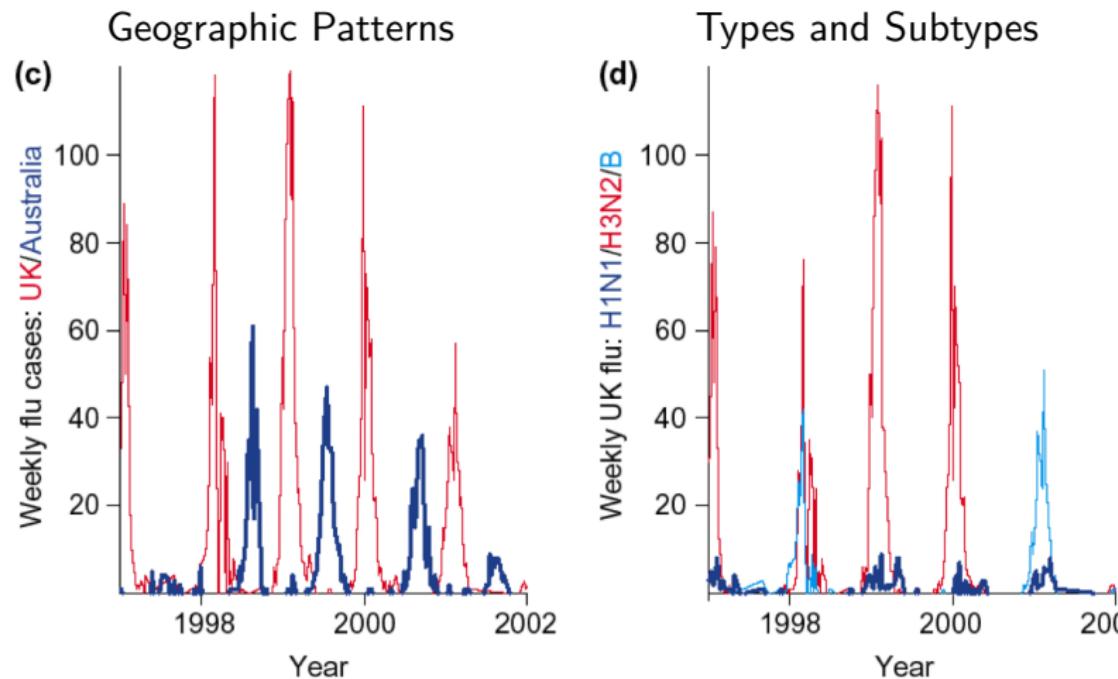
Visualization of entire course of the Great Plague

- What happened after initial spatial spread?
- Visualize full spatial epidemic structure
- Show magnitude of epidemic in each parish with cylinder.
- [Epidemic Visualization](#) (EpiVis) software by Junling Ma.

P&I mortality in U.S.A., 1910–1998



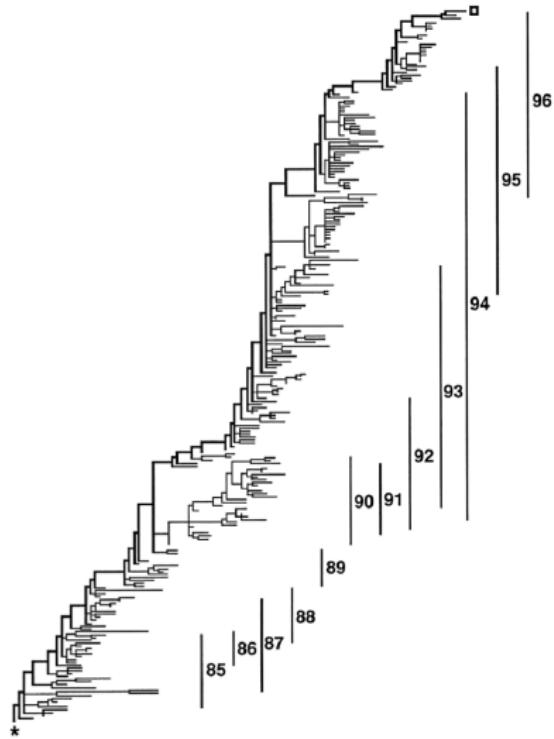
Influenza Incidence Patterns (lab confirmed)



Earn, Dushoff & Levin 2002, *Trends in Ecology and Evolution* 17, 334–340

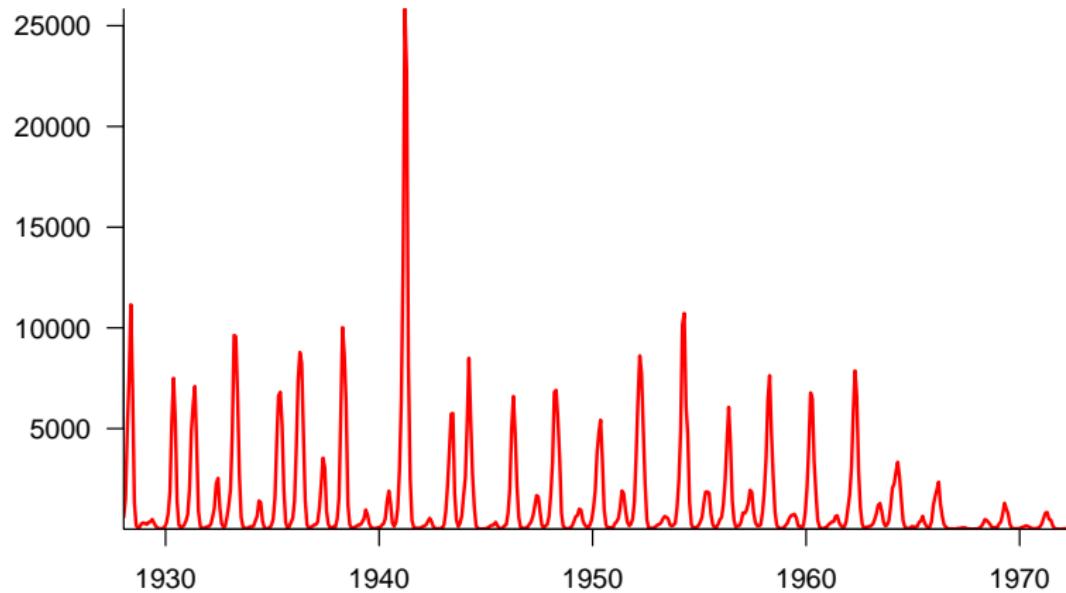
Influenza Evolution

Molecular phylogenetic reconstruction of influenza A/H3N2 evolution, 1985–1996 (Fitch *et al.* 1997)



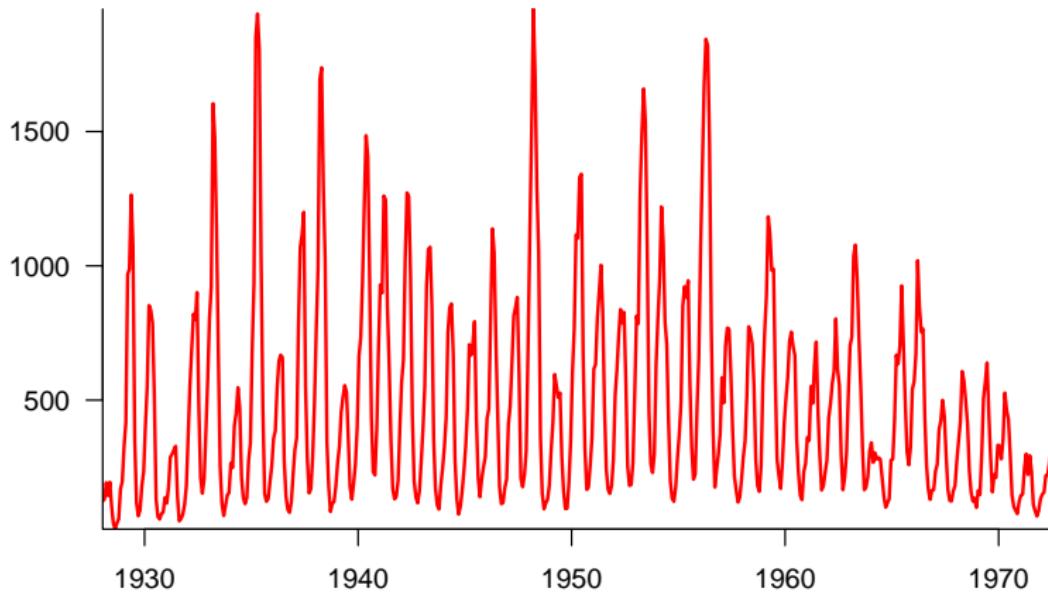
Measles in New York City, 1928–1972

Monthly Cases



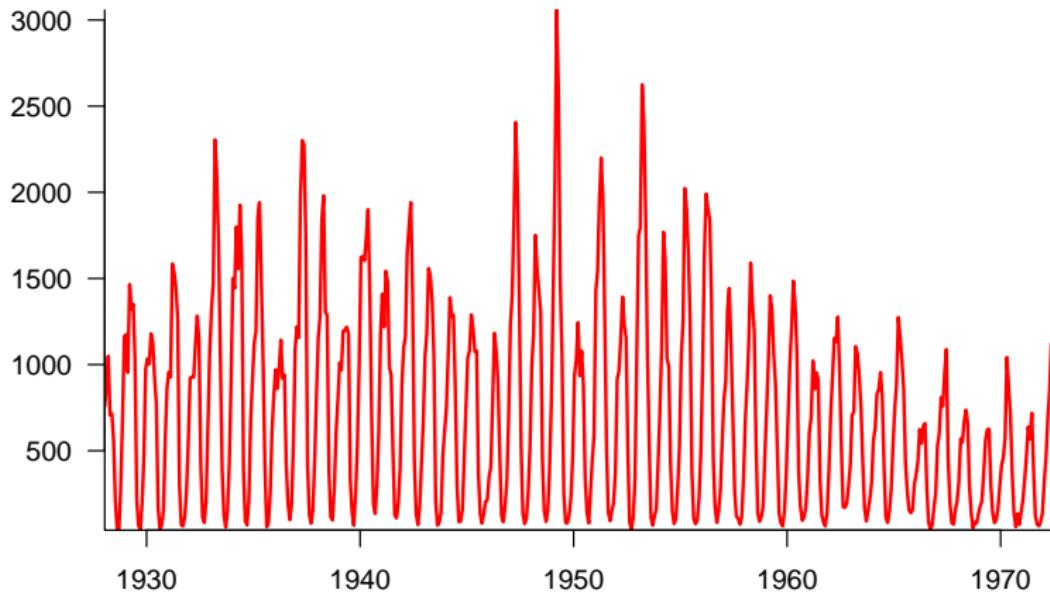
Mumps in New York City, 1928–1972

Monthly Cases

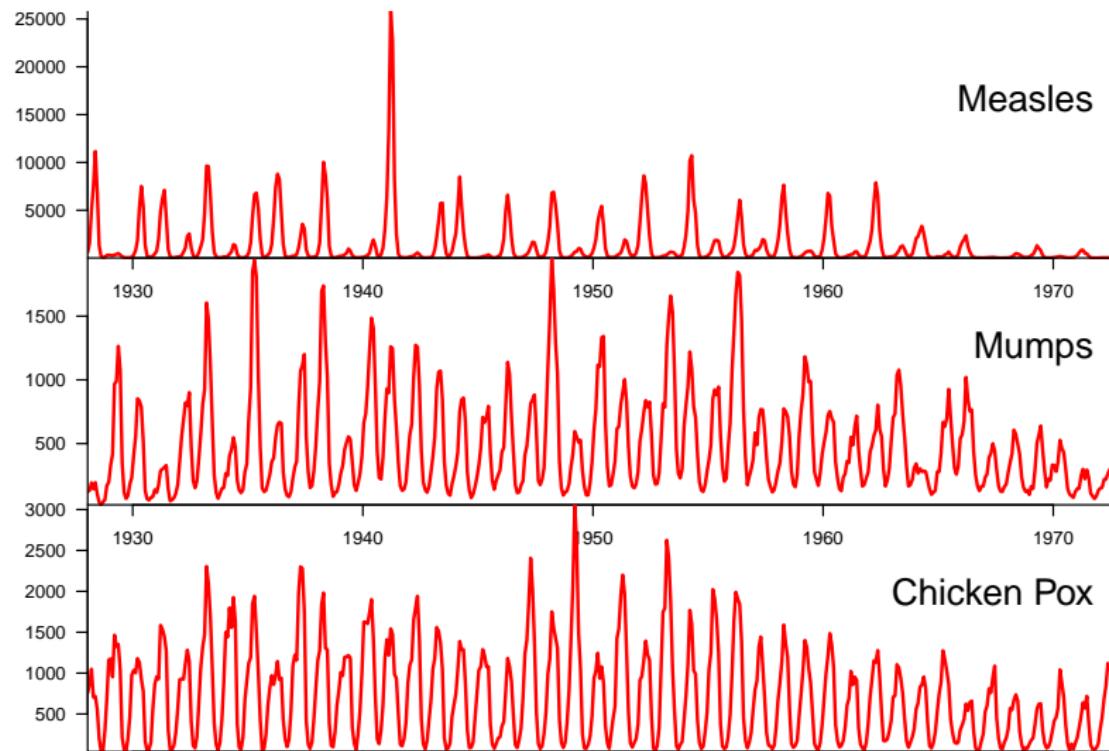


Chicken Pox in New York City, 1928–1972

Monthly Cases

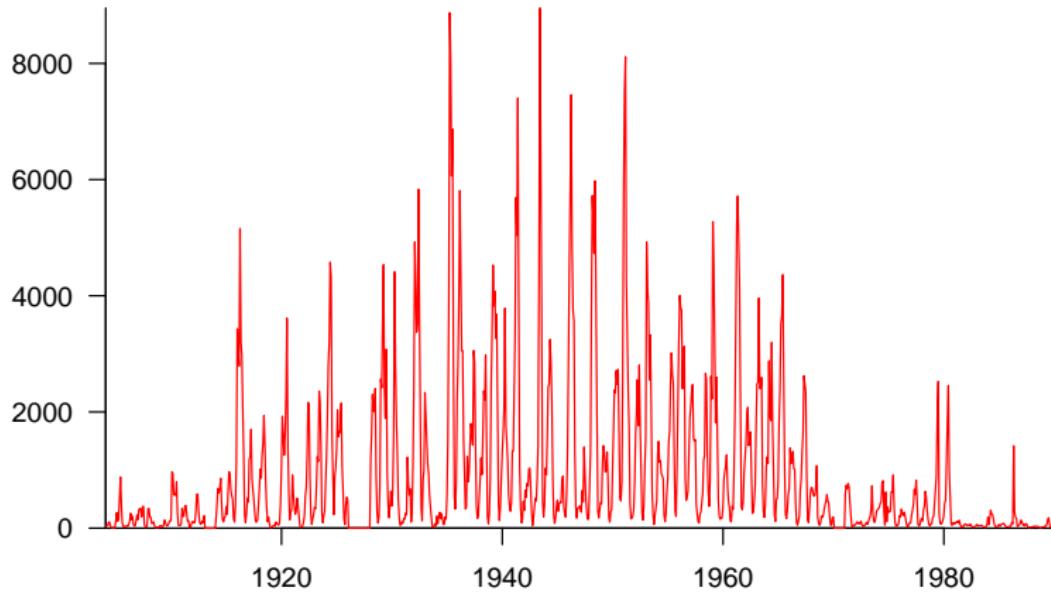


Childhood diseases in New York City, 1928–1972



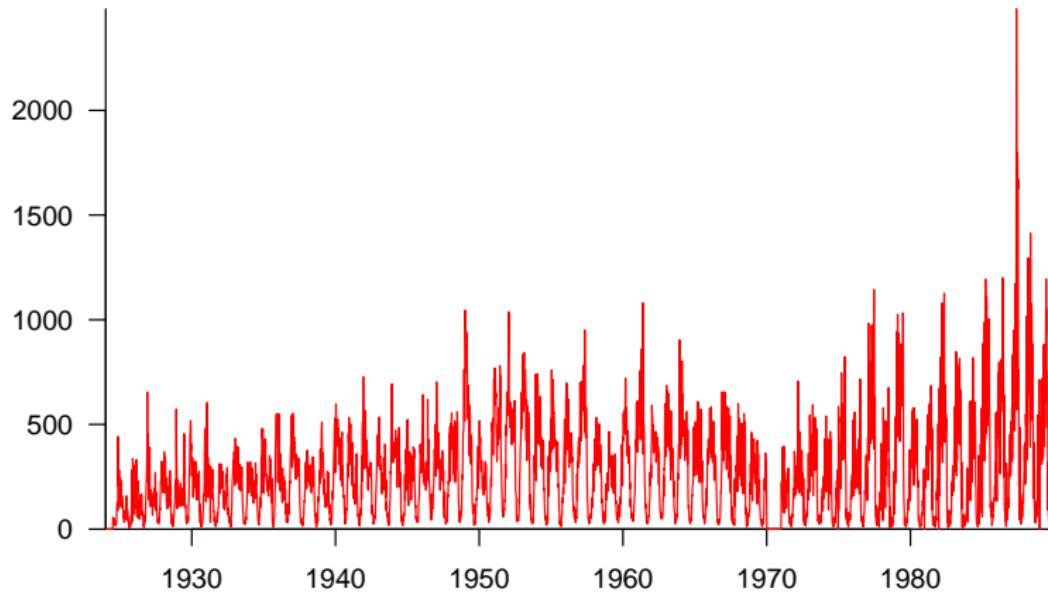
Measles in Ontario, 1904–1989

Monthly Cases



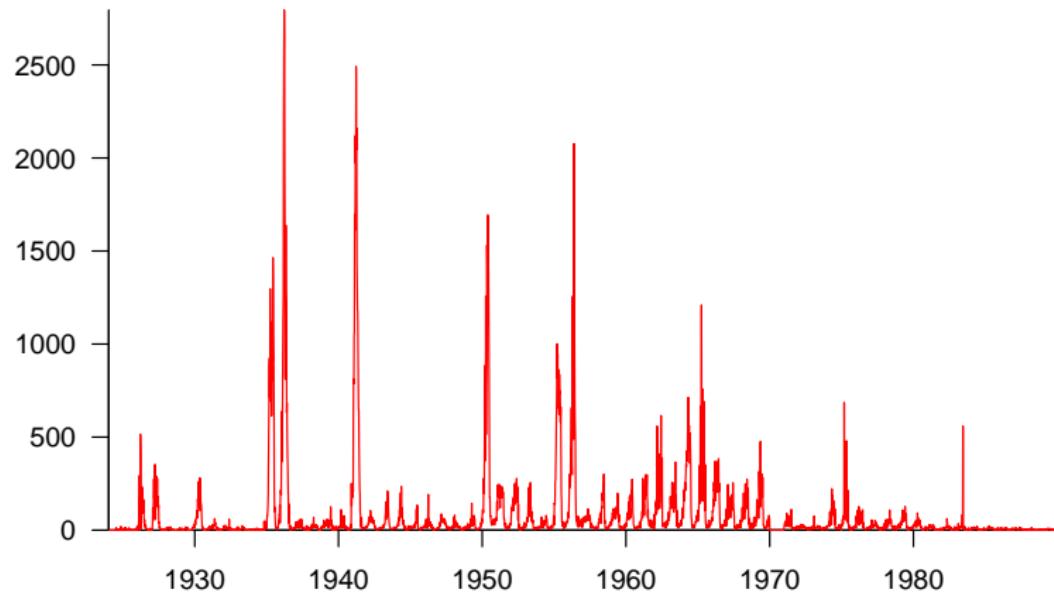
Chicken Pox in Ontario, 1924–1989

Monthly Cases



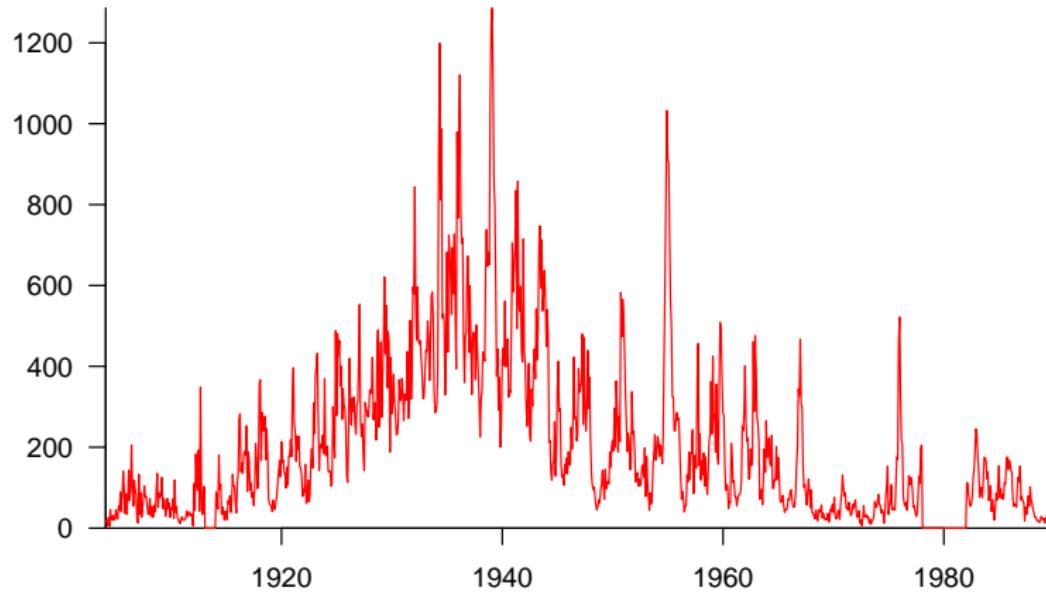
Rubella in Ontario, 1924–1989

Weekly Cases

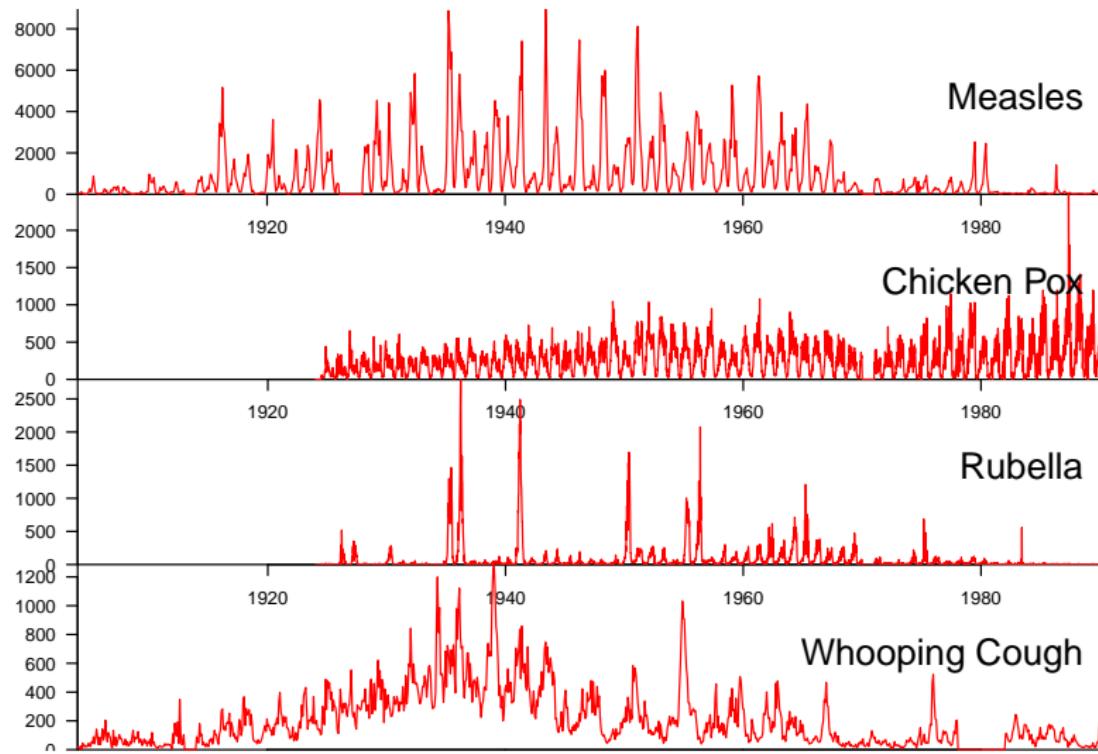


Whooping Cough in Ontario, 1904–1989

Monthly Cases



Childhood diseases in Ontario, 1904–1989



Ontario Disease Notification Data

Dominion Bureau of Statistics Disease Notification Data

VITAL STATISTICS BRANCH - COMMUNICABLE DISEASE SECTION

Cases of Whooping Cough... Reported by Provincial Health Departments, Year 1924

WEEK ENDING	P.E.I.	N.S.	N.B.	QUE.	ONT.	MAN.	SASK.	ALTA.	B.C.	CANADA
WEEK 1	11						1			12
2	12	29					18			47
3	19	37					32			69
4	26	75 152		68	181	36	13 64	97	4 88 602	
5 FEB 2	12	1					53			66
6	1	5					40			45
7	16	31					14			45
8	23	- 2 50 1 2	267	202	48	4 111	116	1	7 797	
9 MAR 1	2						21			23
10	9						9			9
11	15	3					11			14
12	22	60					34			94
13	29	2 61		144	140	52	15 90	15	7 17 515	
14 APR 5	9						11			20
15	12	1					12			13
16	19	26	1				8			35
17	26	14 50 3 4	42	140	39 16 47	67	5	33 394		
18 MAY 3	26						2			28

All Historical Canadian Infectious Disease Data

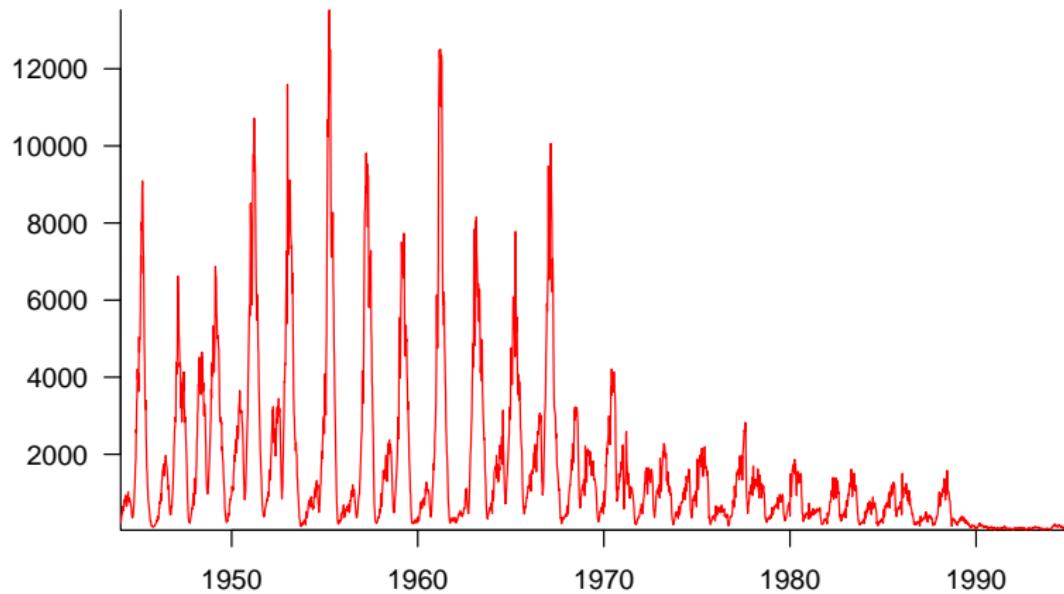
<https://canmod.net/digitization/>

Recurrent epidemics of childhood infections

- Childhood diseases in New York City, 1928–1972
- Childhood diseases in Ontario, 1904–1989

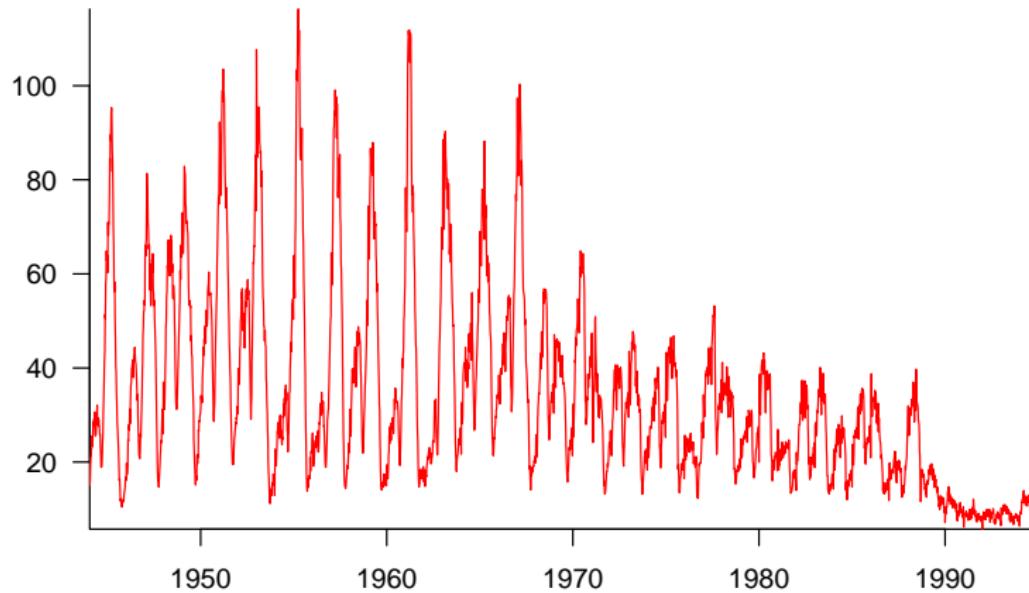
Measles incidence in England and Wales, 1944–1995

Weekly Cases



Measles incidence in England and Wales, 1944–1995

Sqrt(Weekly Cases)



Why study measles epidemics?

- ~ 140,000 annual deaths from measles
- A major cause of *vaccine-preventable* deaths.
- Potential impact in developed countries during vaccine scares (e.g., MMR scare in UK in 1990s).

- Understand past patterns
- Predict future patterns
- Manipulate future patterns
- Develop vaccination strategy that can...



Other reasons to model infectious disease epidemics

- Mathematical models make hypotheses and inferences precise
 - Give better advice to policymakers
 - Make better predictions
- Host-pathogen dynamics are important aspects of ecosystem dynamics
 - Infectious disease models more likely to be successful than predator-prey models
- Excellent data for human infectious diseases
 - Models can be tested!

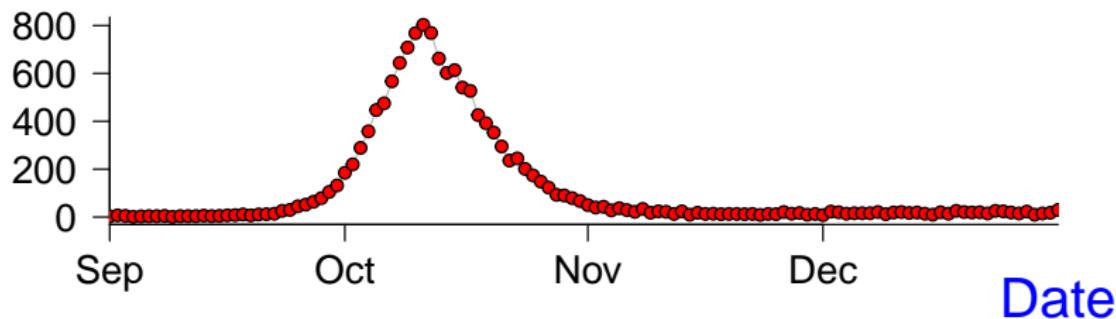
Modelling population dynamics of childhood infections

- The basic SIR model cannot explain recurrent epidemics.
- What should we do?... The usual options:
 - 1 Get depressed, drop the course.
 - 2 Keep developing models until we can explain recurrent epidemics.
- First, let's talk about tools that allow us to make our questions about time series data more precise.

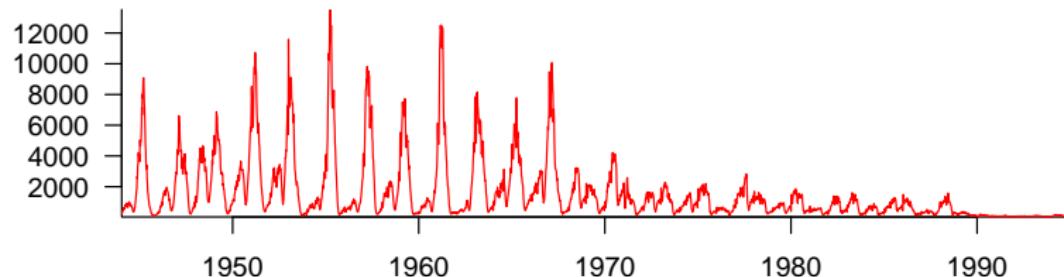
Epidemic Data Analysis

Time Plots of Temporal Epidemic Patterns

1918 P&I

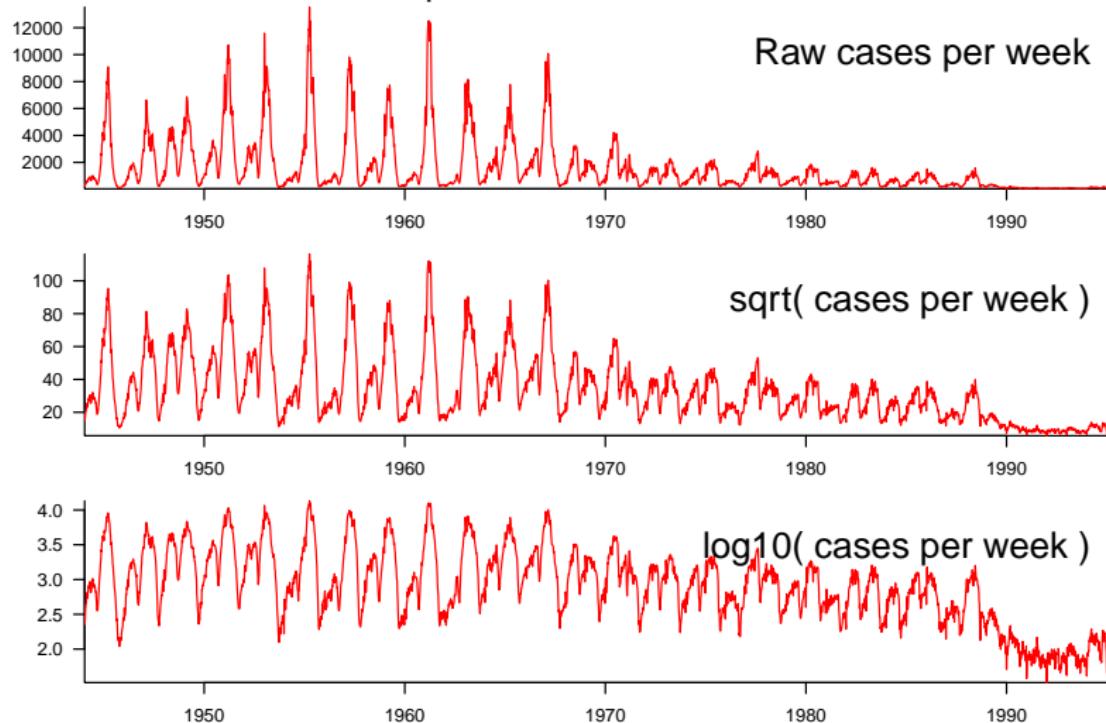


Weekly Measles in England and Wales



Time Plots of Transformed Data

- Reveal unobvious aspects of time series



Times Plots of Smoothed Data

- Reveal trends clouded by noise or seasonality
- *Moving Average:*

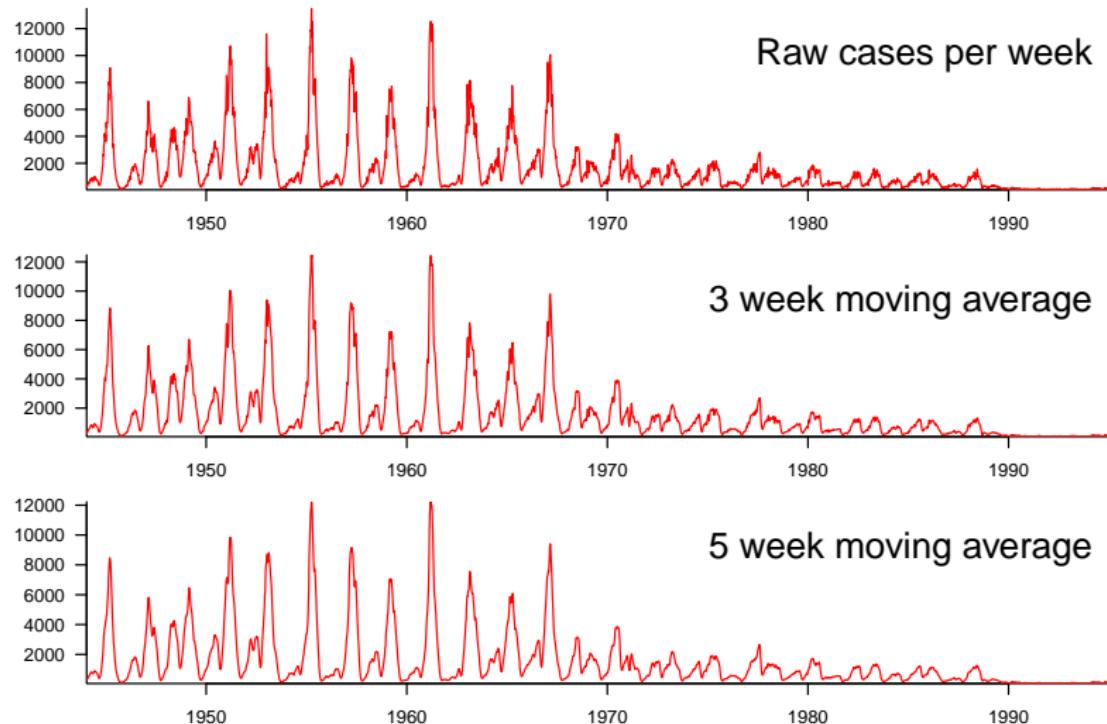
$$x_t \rightarrow \frac{1}{2a+1} \sum_{i=-a}^a x_{t+i}$$

- Replace original data points x_t with averages of nearby points.
- *Linear filter:*

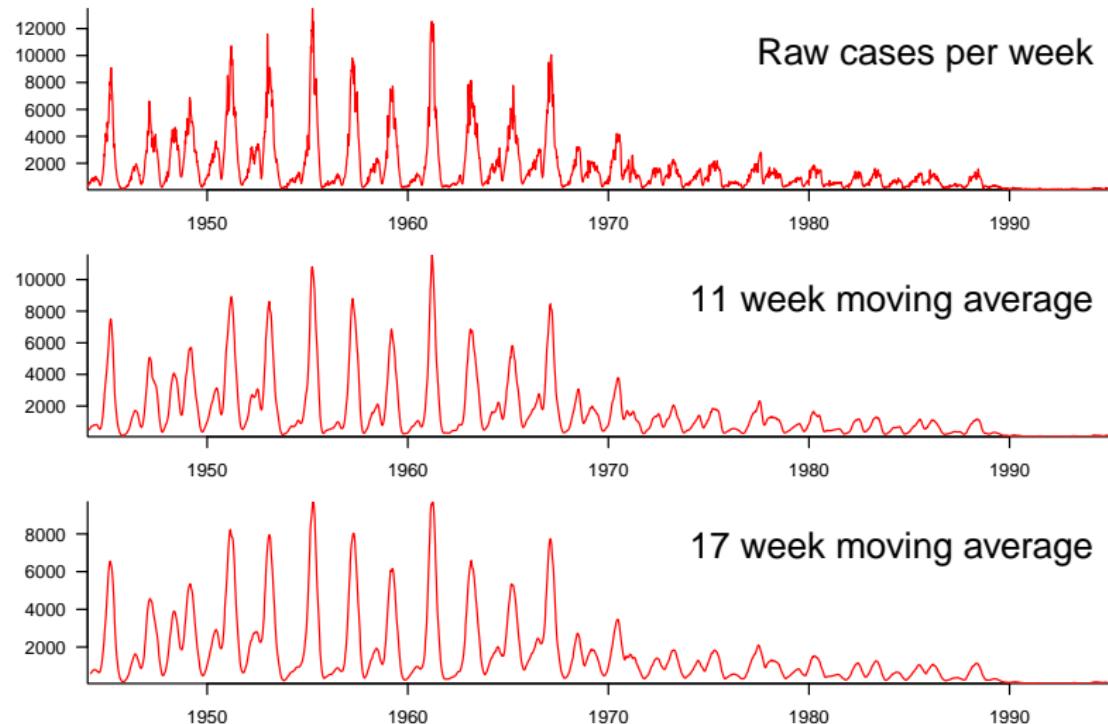
$$x_t \rightarrow \sum_{i=-\infty}^{\infty} \lambda_i x_{t+i}$$

- Generalization of moving average.
- Weights λ_i can be nonlinear functions of i .

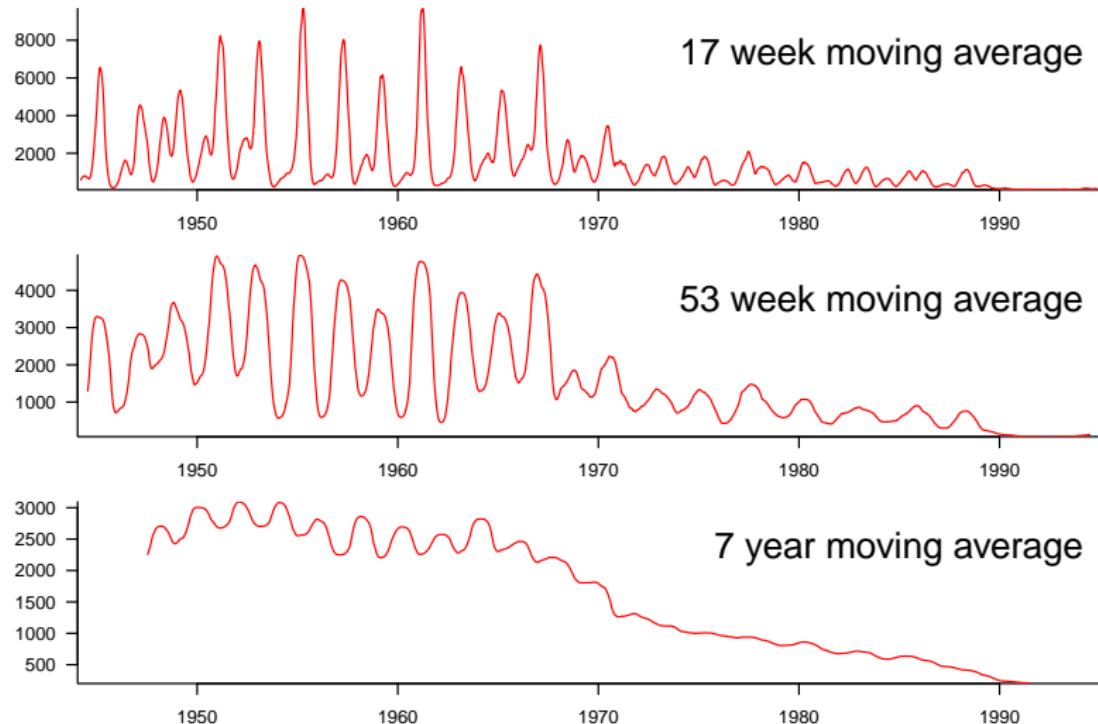
Times Plots of Smoothed Data



Times Plots of Smoothed Data



Times Plots of Smoothed Data



Correlation

- Recurrent epidemics \implies number of cases now is correlated with number of cases in the past and the future.
- Given N pairs of observations of different quantities, $\{(x_i, y_i) : i = 1, \dots, N\}$, the *correlation coefficient* is defined to be

$$r = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2 \sum_{i=1}^N (y_i - \bar{y})^2}}$$

where \bar{x} and \bar{y} are the means of $\{x_i\}$ and $\{y_i\}$, respectively.

Correlation

Properties of the correlation coefficient:

- $-1 \leq r \leq 1$ (Proof? Cauchy-Schwarz inequality)
- $r = 1 \iff$ all points lie on a line with positive slope ("complete positive correlation")
- $r = -1 \iff$ all points lie on a line with negative slope ("complete negative correlation")
- $r \simeq 0 \implies$ "uncorrelated"
- *Interpretation:* r^2 is the proportion of the variance in y explained by a linear function of x .

Derivations and discussions:

- [MathWorld on \$r^2\$](#) , [Wikipedia on \$r^2\$](#)
- [Wikipedia on general coefficient of determination](#)

Autocorrelation

- Given a single sequence of observations $\{x_t : t = 1, \dots, N\}$, we can compute the correlation of each observation with the observation k time steps in the future.
- Thus, we consider the pairs of observations $\{(x_t, x_{k+t}) : t = 1, \dots, N - k\}$ and define the *autocorrelation coefficient at lag k* to be

$$r_k = \frac{\sum_{t=1}^{N-k} (x_t - \bar{x}_{1,N-k})(x_{k+t} - \bar{x}_{k+1,N})}{\sqrt{\sum_{t=1}^{N-k} (x_t - \bar{x}_{1,N-k})^2 \sum_{t=1}^{N-k} (x_{k+t} - \bar{x}_{k+1,N})^2}}$$

where $\bar{x}_{1,N-k}$ and $\bar{x}_{k+1,N}$ are the means of first and last $N - k$ observations, respectively.

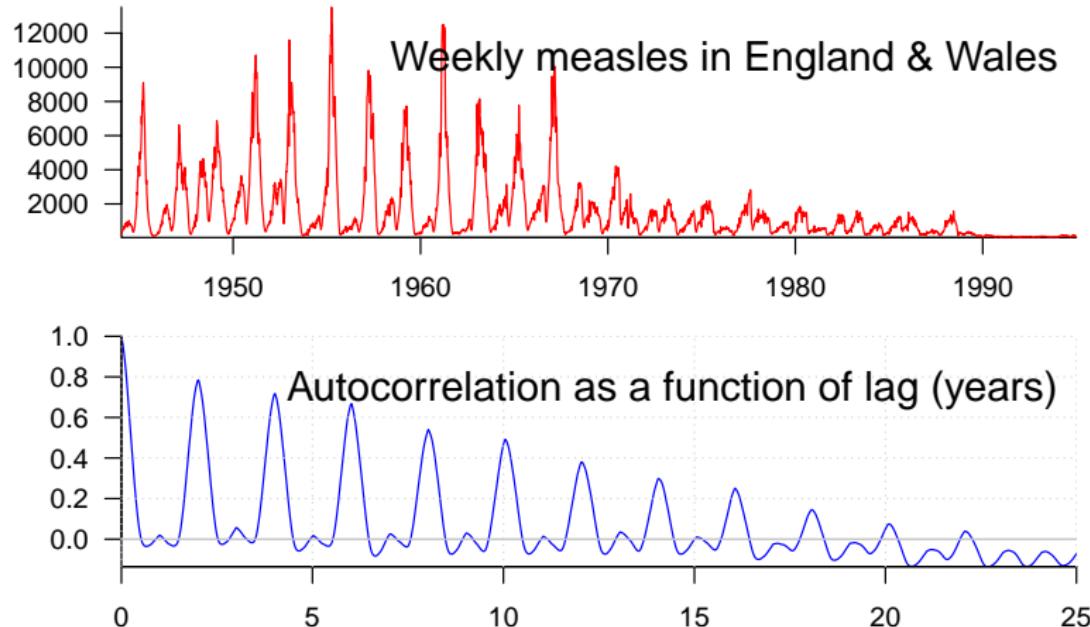
Autocorrelation

- If number of observations N is large and lag $k \ll N$ then

$$r_k \simeq \frac{\sum_{t=1}^{N-k} (x_t - \bar{x})(x_{k+t} - \bar{x})}{\sum_{t=1}^N (x_t - \bar{x})^2}$$

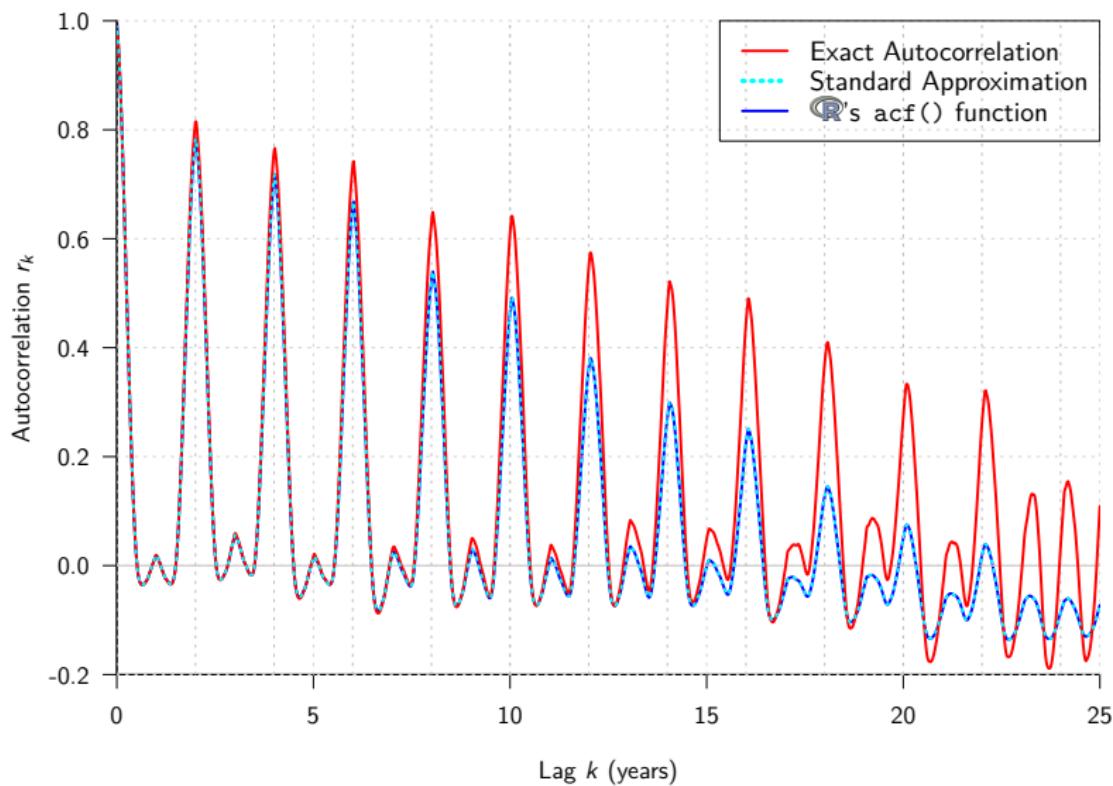
- Approximation of r_k is worse for larger lags k
- Plot of autocorrelation r_k as a function of lag k is called the *correlogram*.

Correlogram



- Peaks in correlogram \implies periodicities in original time series.
- Correlograms of temporal segments are often informative.

Correlogram: exact vs. approximate r_k



Spectral Density

- Can we compute the dominant periods in the time series?
(Rather than estimating them by eye from the [correlogram](#).)
- Express the time series as a [Fourier series](#):

$$x_t = a_0 + \left(\sum_{p=1}^{(N/2)-1} (a_p \cos \omega_p t + b_p \sin \omega_p t) \right) + a_{N/2} \cos \pi t,$$

where $\omega_p = 2\pi p/N$.

- Compute the [Fourier coefficients](#) $\{a_p\}$, $\{b_p\}$ by taking inner products with $\cos \omega_p t$ and $\sin \omega_p t$.

Spectral Density

- Fourier coefficients of x_t are:

$$a_0 = \bar{x} = \frac{1}{N} \sum_t x_t ,$$

$$a_p = \frac{2}{N} \sum_t x_t \cos \omega_p t , \quad b_p = \frac{2}{N} \sum_t x_t \sin \omega_p t ,$$

$$a_{N/2} = \frac{1}{N} \sum_t (-1)^t x_t ,$$

where sum is over observation times.

- Estimated power spectral density (PSD) at frequency ω_p is*:

$$I(\omega_p) = \frac{N}{4\pi} (a_p^2 + b_p^2)$$

*The normalization by $N/4\pi$ is the convention chosen by Chatfield (2004, "Analysis of Time Series: An Introduction"). Other normalization conventions are also in common use.

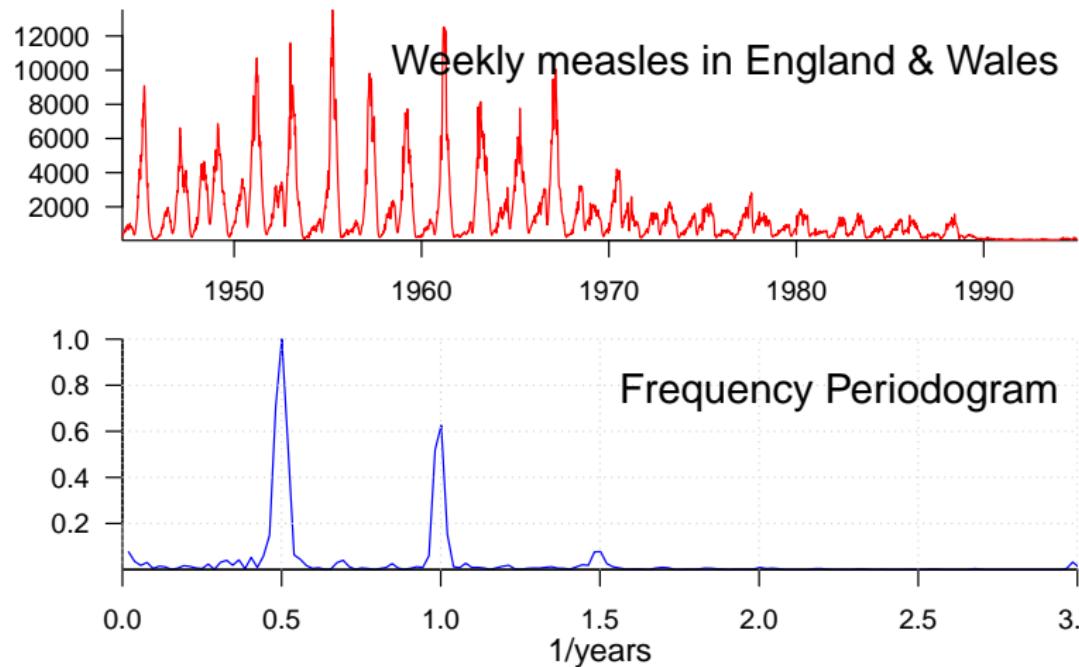
Spectral Density

- There are many different ways to express the power spectral density (aka *power spectrum*).
- Most common/useful equivalence is that the power spectrum is the discrete Fourier transform of the correlogram:

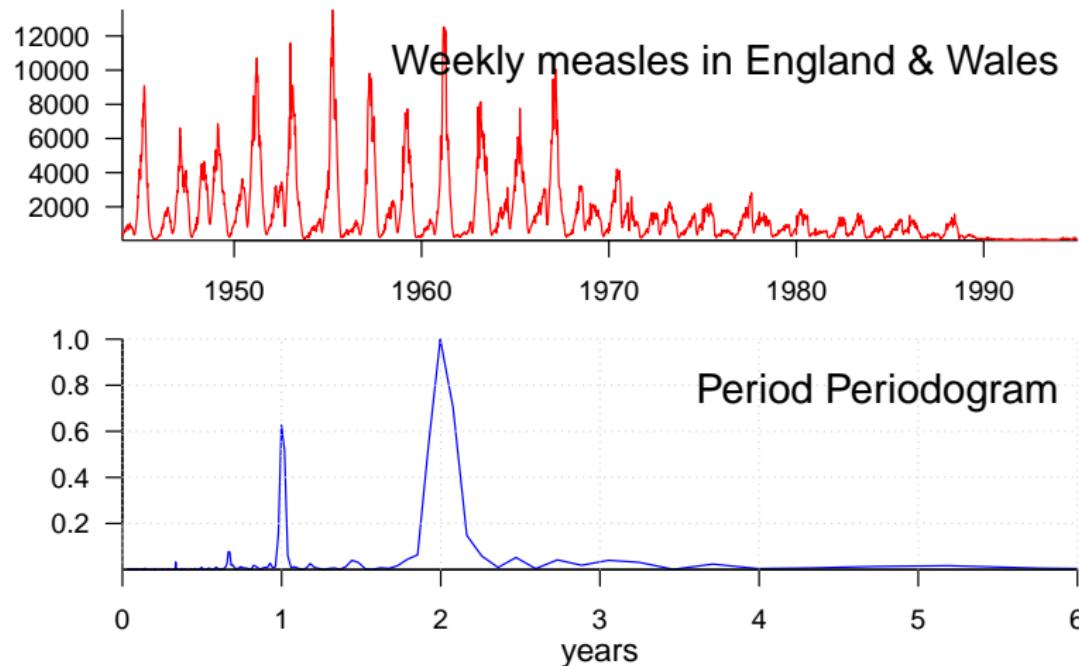
$$I(\omega_p) = \frac{1}{\pi} \left(r_0 + 2 \sum_{k=1}^{N-1} r_k \cos \omega_p k \right)$$

- Plot of estimated power spectrum as a function of frequency ω_p is called the *frequency periodogram* or just the *periodogram*.

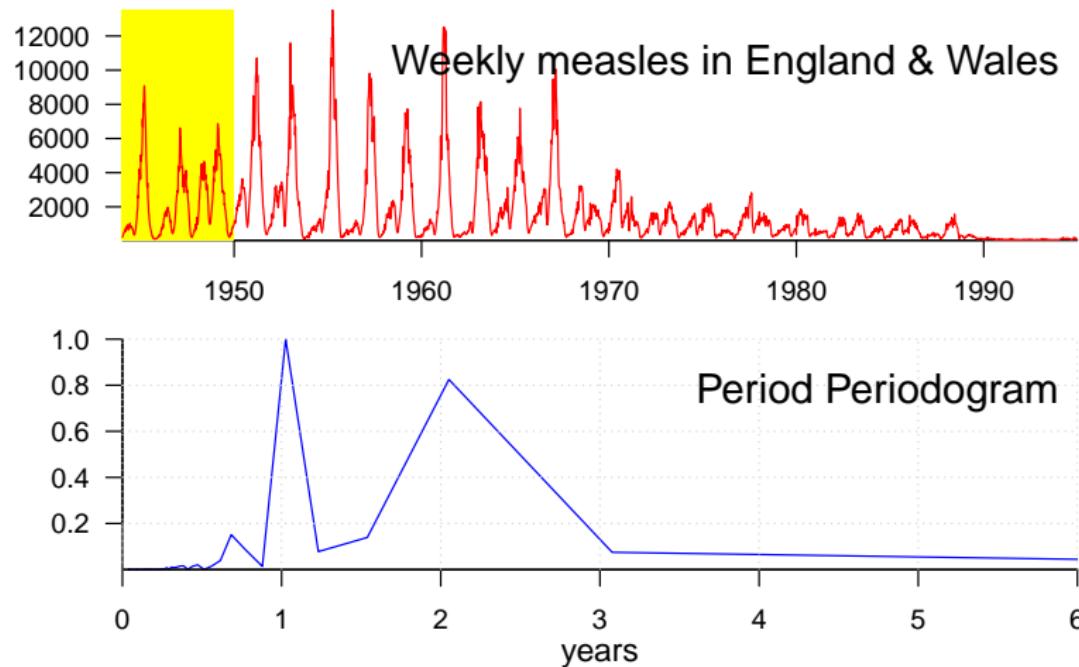
Spectral Density



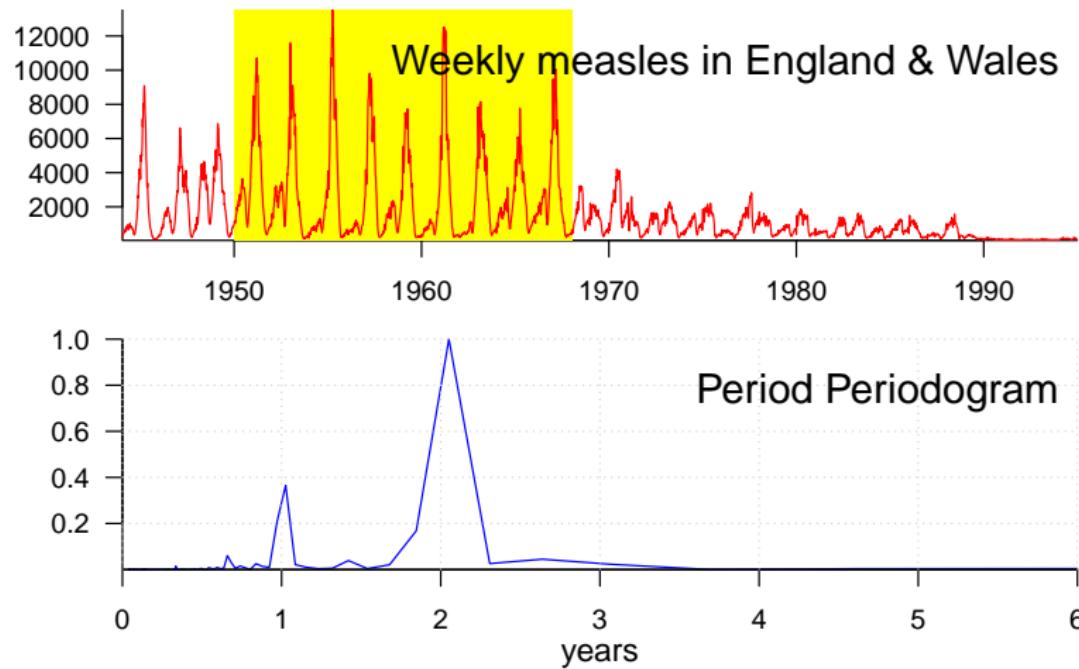
Spectral Density



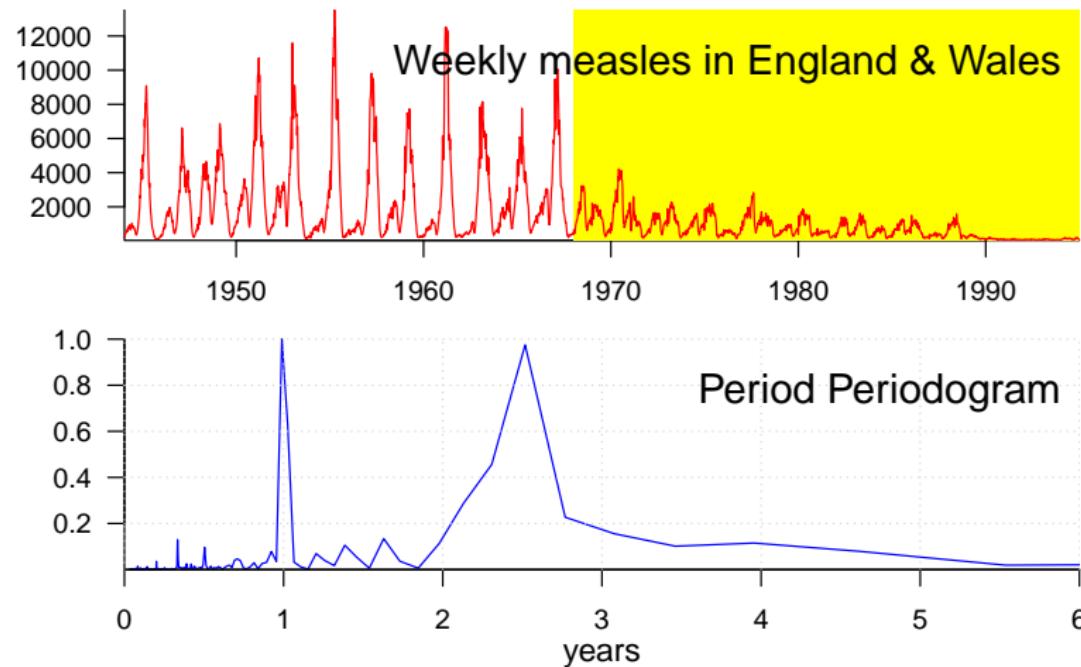
Spectral Density of Temporal Segments



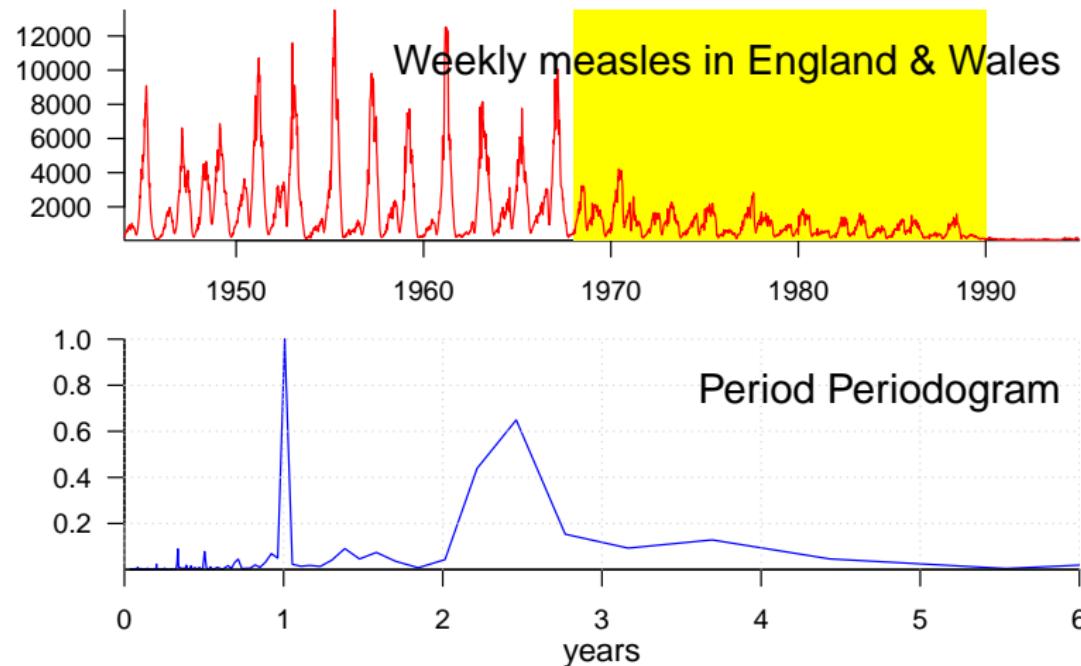
Spectral Density of Temporal Segments



Spectral Density of Temporal Segments



Spectral Density of Temporal Segments

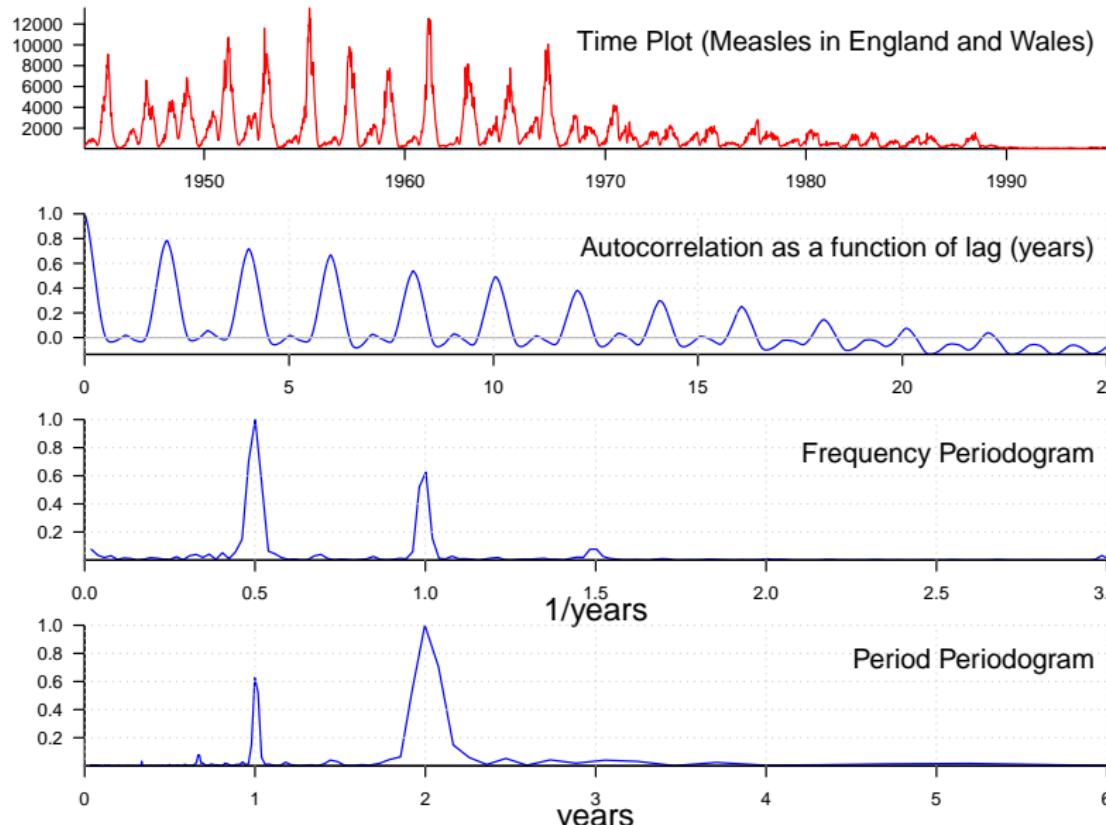


Spectral Density Properties

- Periodogram is discrete Fourier transform of correlogram
- Same information in correlogram and periodogram
- Periodogram usually easier to interpret
- In , calculate power spectrum with `spectrum()`
- The power spectrum $I(\omega_p)$ partitions the variance in the time series with respect to frequency ω_p .
 - Parseval's theorem implies $\frac{1}{N} \sum_t (x_t - \bar{x})^2 = \frac{1}{2\pi N} \sum_{p>0} I(\omega_p)$.
But $\frac{1}{N} \sum_t (x_t - \bar{x})^2 = \text{Var}\{x_t\}$, hence $I(\omega_p)/(2\pi N)$ is the proportion of the variance in the time series associated with period $2\pi/\omega_p$.

[For details, see Chatfield (2004).]

Basic Time Series Analysis of Epidemic Data



Announcements

Spectral Density of Temporal Segments

- Pre-war measles
- Post-war pre-vaccination measles
- Vaccination era measles
- Vaccination era measles until 1990

Time series analysis functions



has built-in tools for time series analysis:

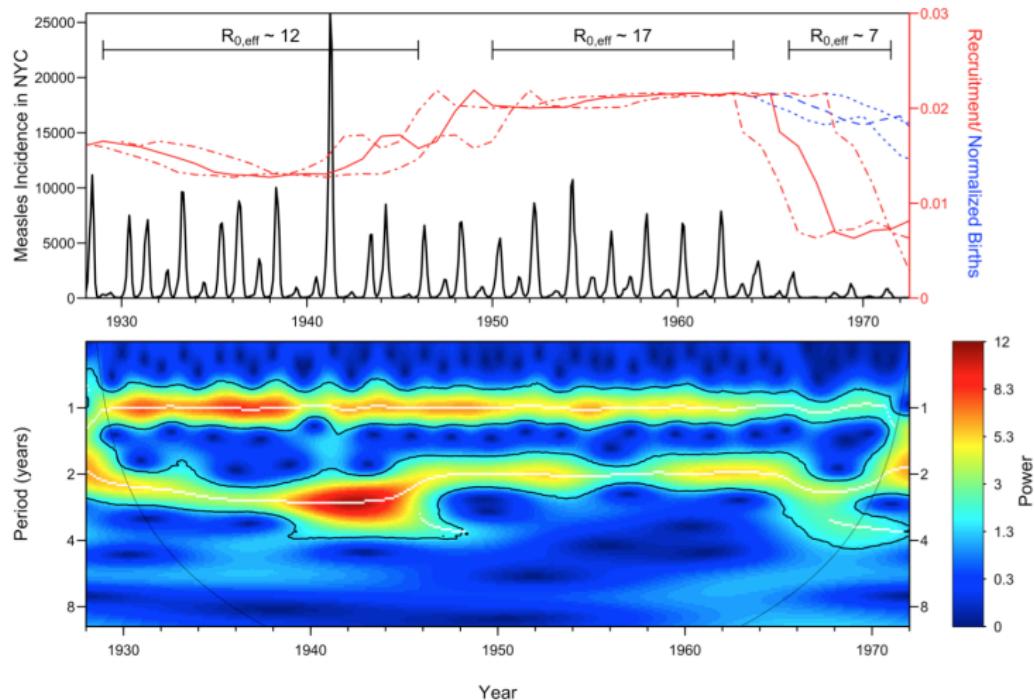
- Time plot: `plot()` etc.
- Linear filter (e.g., moving average): `filter()`
- Correlogram (auto-correlation function): `acf()`
- Periodogram (power spectrum): `spectrum()`

You will use all of these functions in **Assignment 4**.

More sophisticated spectral method

- Traditional power spectrum measures frequency content of entire time series.
- Wavelet decomposition is local in time.
 - Reveals changes in the spectrum over time without having to identify distinct temporal segments yourself.
 - Nice intro to wavelet analysis of time series:
Torrence and Compo (1998) "A Practical Guide to Wavelet Analysis" *Bulletin of the American Meteorological Society* **79**, 61–78
 - $\exists \text{ } \text{R}$ packages for wavelet analysis of time series (e.g., `WaveletComp`, `wavelets`), and at least one book on wavelet methods in 

Wavelet Spectrum of Monthly Measles in New York City



Krylova & Earn 2013, *J. R. Soc. Interface* **10**, 20130098

Wavelet Spectrum of Weekly Measles in New York City

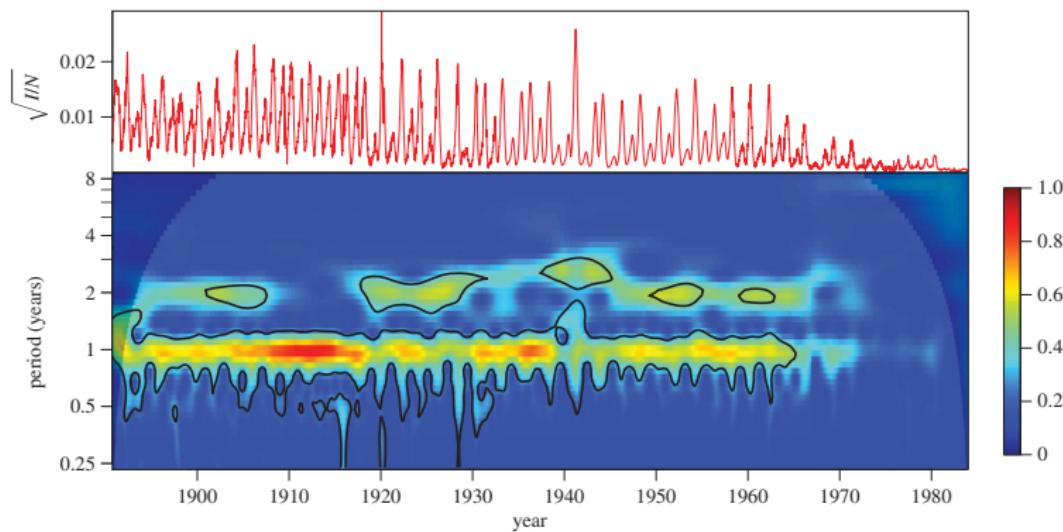
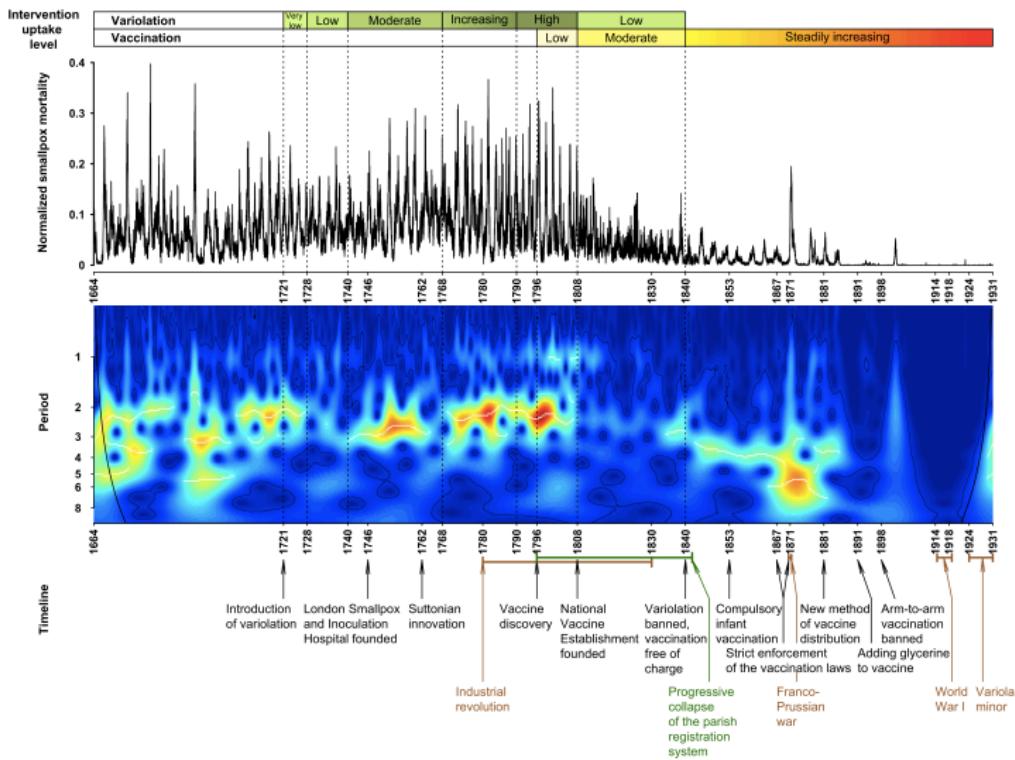


Figure 5. Observed measles dynamics in NYC from 1891 to 1984. (a) Square root of measles case reports, normalized by total concurrent population. (b) Colour depth plot of a continuous wavelet transform of the square root of normalized observed NYC measles cases (colour warmth scales with spectral power and 95% significance contours are shown in black). Shaded regions in the upper left and right indicate the cone of influence.

Hempel & Earn 2015, *J. R. Soc. Interface* 12, 20150024

Wavelet Spectrum of Weekly Smallpox in London



Krylova & Earn 2019, *bioRxiv* doi: <https://doi.org/10.1101/771220>

Statistical Modelling of Time Series

Statistical Modelling of Time Series

- Imagine time series $\{X_t\}$ is generated by random processes.
- Simplest case: X_t (number of cases at time t) is simply a random variable with a known distribution,

$$X_t = \mu + Z_t \quad (*)$$

where μ = time average number of cases
and $\{Z_t\}$ = sequence of random variables with zero mean.

- Might be a reasonable model for importation of new, infectious individuals into a focal community.
- Bad model for epidemics: ignores transmission from one individual to another.
 - There must be a correlation between the number of individuals in the focal community who are infected now and the number who will be infected in the near future.

Statistical Modelling of Time Series: AR and MA

- So, imagine that successive data points in $\{X_t\}$ are correlated.
- For example, perhaps the data are generated by an *autoregressive (AR) process*:

$$X_t - \mu = \alpha_1(X_{t-1} - \mu) + \alpha_2(X_{t-2} - \mu) + \cdots + \alpha_p(X_{t-p} - \mu) + Z_t,$$

where the α_i are constants that determine the degree of correlation along the time series.

- Alternatively, the data might be generated by a *moving average (MA) process*:

$$X_t - \mu = \beta_0 Z_t + \beta_1 Z_{t-1} + \cdots + \beta_q Z_{t-q},$$

where the β_i are constants that define a weighted average.

Statistical Modelling of Time Series: ARMA

- More generally, the data might be generated by an *autoregressive moving average “ARMA(p, q)” process:*

$$\begin{aligned} X_t - \mu = & \alpha_1(X_{t-1} - \mu) + \alpha_2(X_{t-2} - \mu) + \cdots + \alpha_p(X_{t-p} - \mu) \\ & + \beta_0 Z_t + \beta_1 Z_{t-1} + \cdots + \beta_q Z_{t-q}. \end{aligned}$$

Statistical Modelling of Time Series: ARIMA

- Finally, an *autoregressive integrated moving average “ARIMA(p, d, q)” model* includes weighted differences of the time series:

$$\begin{aligned} X_t - \mu &= \alpha_1(X_{t-1} - \mu) + \alpha_2(X_{t-2} - \mu) + \cdots + \alpha_p(X_{t-p} - \mu) \\ &\quad + \gamma_1(X_{t-1} - X_{t-2}) + \gamma_2(X_{t-2} - X_{t-3}) + \cdots \\ &\quad + \beta_0 Z_t + \beta_1 Z_{t-1} + \cdots + \beta_q Z_{t-q}. \end{aligned}$$

- The “I” in ARIMA refers to the original time series X_t , which is an “integrated” version of the differenced time series.
- Technically, an ARIMA model is just an ARMA model with differently labelled coefficients, but explicit differences are often helpful conceptually (e.g., they can “stationarize” a time series).

What kind of process generated our data?

- *How can we tell if our data were generated by such a process?
Can we identify an AR(p), MA(q) or ARMA(p, q) process?*

- Compare time plots of these processes with time plot of our data? (Comparison by eye often challenging/unreliable.)
- Compare autocorrelation functions (correlograms) of these processes with correlogram of our data? (Better.)
- Compare power spectra (periodograms) of these processes with periodogram of our data? (Even better.)
- Compare wavelet spectra of these processes with wavelet spectrum of our data? (Better yet.)

Statistical Modelling of Time Series: ARMA fitting

- Looking at the power spectra of ARMA models would be instructive.
- But is there a better approach to discovering if an ARMA model could explain our data?
- Find the *best fit* ARMA parameters by minimizing the residual sum of squares. e.g., for an AR model, minimize:

$$S = \sum_{t=p+1}^N [(x_t - \mu) - \alpha_1(x_{t-1} - \mu) - \cdots - \alpha_p(x_{t-p} - \mu)]^2.$$

- More generally, we can find the best fit parameters of an ARIMA(p, d, q) model
 - Non-trivial, but there are standard methods
- Compare models with **Akaike Information Criterion (AIC)**, which penalizes models that have more parameters
 - See [Earn \(2009\)](#) review article for more discussion of this.

Time series tools discussed so far...

- Statistical description of time series:
time plot, moving average, correlation coefficient,
autocorrelation, correlogram, power spectral density (PSD),
periodogram, wavelet spectrum
- Time series models:
AR, MA, ARMA, ARIMA

Statistical Modelling of Time Series

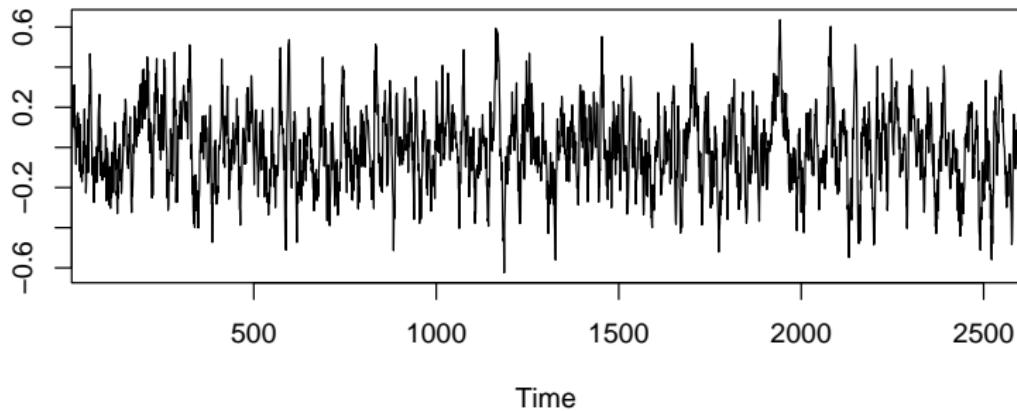
How to do it in ...

- Simulate any ARIMA(p, d, q) model with `arima.sim()`
- Fit an AR model to a time series with `ar()`
- Fit an ARIMA model to a time series with `arima()`
- Alternatively, there are specialized time series modelling packages.

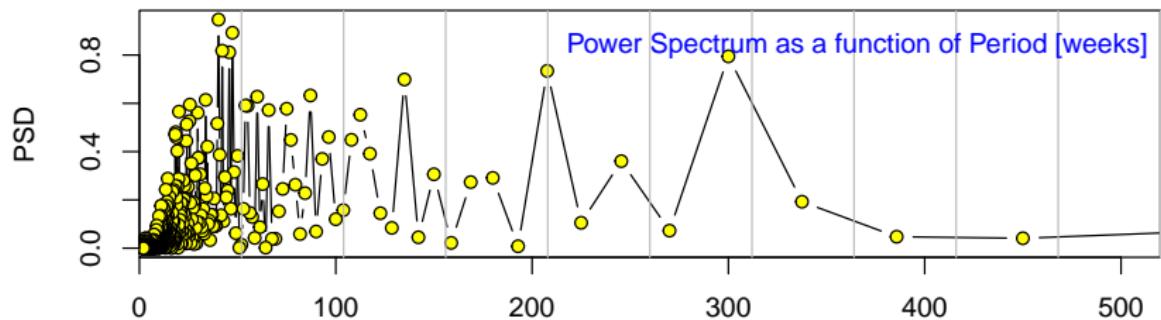
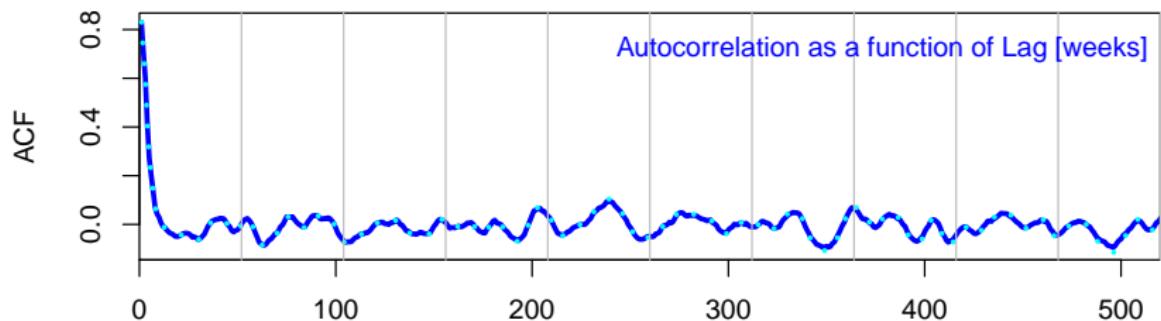
ARMA Example (50 years of weekly data)

```
my.model <- list(ar=c(1,-0.5,0.5,-0.25),ma=c(-0.25,0.5))
my.sim <- arima.sim(n=52*50,model=my.model,sd=0.1)
plot(my.sim,main="ARMA Example",ylab="",xaxs="i")
```

ARMA Example



ARMA Example (ACF and PSD up to 10 year lag)



Statistical Modelling of Time Series: Forecasting

- Once we have a fitted model, we can then use it to *forecast* future observations
- *Validate* this procedure by using part of the data to fit the model and then forecast the remainder of the data (*cf.* cross-validation)
- How successful is this likely to be for an infectious disease time series?
 - Conceivably good for chicken pox in NYC.
 - Less likely to be good for measles... at least for the main patterns...
 - One of the project options is to look at this more carefully.

Statistical Modelling of Time Series: Limitations

- It might be best to remove mean, trend and seasonality before fitting an ARMA model
 - But this means we will remove the aspects of the data about which we care most!
- The fitted parameters of an ARMA model have no obvious biological meaning
 - The model completely ignores any understanding we have of infectious disease transmission
- Statistical models use the time series itself to parameterize an ARMA (or more general) process
 - It would be better to have a model that we can parameterize from independently collected data and then see if that model can explain the observed time series

Mechanistic Mathematical Modelling

- SIR and all that...
- Takes into account transmission process...
- So why did we just spend time talking about statistical modelling?
 - Important to be familiar with time series models that are in common use.
 - Helps us appreciate the value of mechanistic modelling.
 - Some processes that affect disease dynamics might be better modelled as ARMA or similar processes.
 - Weather (e.g., perhaps model $\beta = \beta(t)$ as an ARMA process)
 - Immigration
 - Ruling out an ARMA model (or at least one with a modest number of parameters) is a step towards finding a good model.
 - A combination of mechanistic and time series models could be useful.