

Time-series regression for short-term health effects of environmental exposures in the EMME

Basic concepts for the time-series design

Aurelio Tobías

Institute of Environmental Assessment and Water Research (IDAEA),
Spanish Council for Scientific Research (CSIC)

Pre-conference workshop for the **2nd International Conference on Climate Change in the Eastern Mediterranean and Middle East**

Cyprus/online – 11th October 2021

Outline

- Time-series design
 - Principles of time-series regression
 - Exposure-response and lagged effects
 - Environmental risk exposures
-
- Workshop materials available at:
<https://github.com/aureliotobias/emme2021>

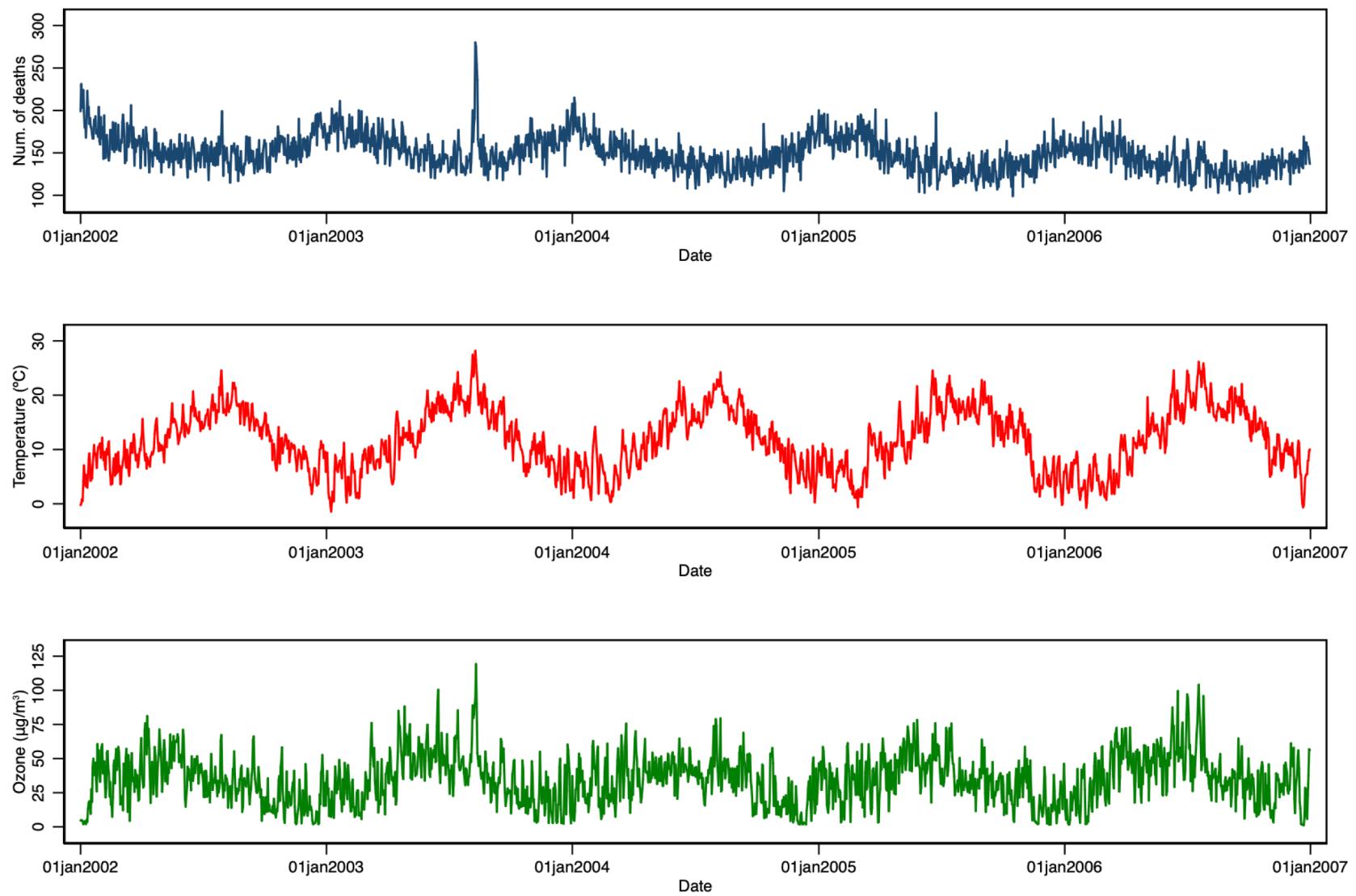
Introduction

- Health effects are the changes in health status resulting from exposure to a given risk factor
 - **Short term effects** – Acute impact on health after an immediate exposure (time-series studies)
 - **Long term effects** – Chronic health effect after a cumulative exposure (cohort studies)
- **Health impact assessment** is the evaluation of potential health effects of proposed actions relative to a given exposure. The aim of HIA is to provide recommendations for decision-making process that will protect health

Introduction

- **Research question** – “*Is there an association between day-to-day variation in the environmental exposure and daily risk of health outcome*”?
 - Health outcomes and environmental exposures are characterized by **similar time-trends**
 - Measures of **individual predictors** are usually not available
 - We need a study design that relies on **between-day comparison within the same population** and able to control for time-trends
 - **Time-series data should be long enough** to identify day-to-day variation necessary to disentangle short-term effects from time-trends (e.g., at least 3 consecutive years for daily data)

London, Jan 2002-Dec 2006



Brief history

- 1952 – Air pollution episode (fog) in London
- 1990's – Time-series regression in air pollution multi-centre studies in EU (APHEA) and US (NMMAPS)
(Katsouyanni et al. 1996, Samet et al. 2000)
- 2000's – Generalised Additive Models (*Hasite and Tibshirani 1990*) in **Splus** (**gam** function)
- 2003 – Heatwave in Europe
- 2010 – Distributed lag non-linear models (*Gasparrini et al. 2010*) in **R** (**dlnm** library)
- 2021 – Case time-series design (*Gasparrini 2021*)

Time-series data

- A time-series is a sequence of measurements equally spaced through time (Zeger *et al.* 2006)
- The unit of analysis used to be the day (t), not the individual person (i)
- But it could be annual, monthly, weekly or hourly
- The outcome is a count (e.g., number of deaths)
- Example – First 10 rows of time-series data (London, Jan 2002 – Dec 2006)

obs	date	deaths	temp	ozone
1.	01jan2002	199	-0.2	4.6
2.	02jan2002	231	0.1	4.9
3.	03jan2002	210	0.9	4.7
4.	04jan2002	203	0.5	4.1
5.	05jan2002	224	4.2	2.0
6.	06jan2002	198	7.1	2.4
7.	07jan2002	180	5.2	4.1
8.	08jan2002	188	3.5	3.1
9.	09jan2002	168	3.2	2.1
10.	10jan2002	194	5.3	5.2

Time-series design

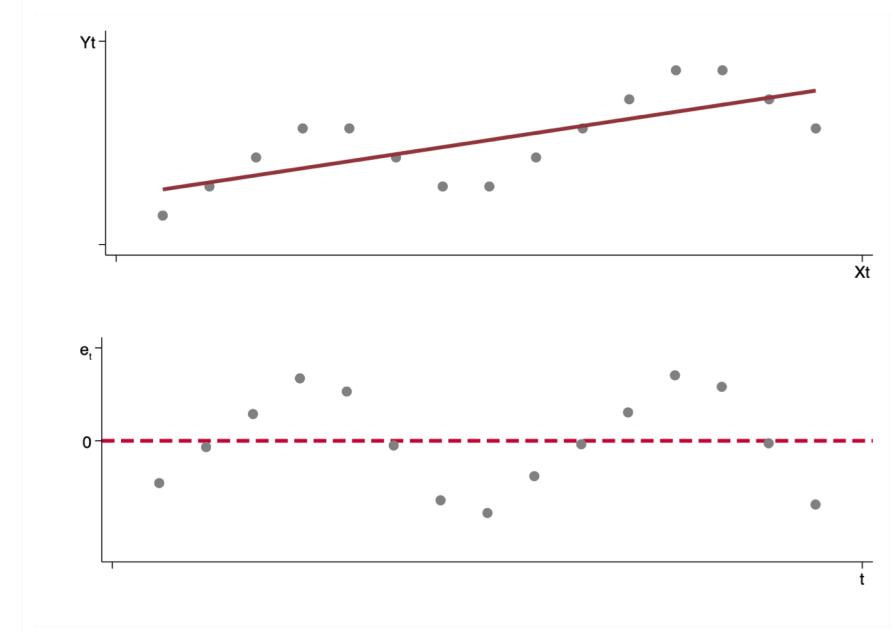
- **Strengths**
 - Use of administratively collected data
 - Same population is compared with itself, focus is day-to-day variation
 - Time-invariant or slowly varying individual risk factors controlled by design (e.g., age, gender, smoking)
- **Limitations**
 - Ecological design based on aggregated, not individual data
 - Not applicable to estimate long-term (chronic) effects
 - Sensitive to choices for modelling time-trends

Poisson regression

- Similar in principle to any regression analysis but with **some specific features**
- Poisson regression
 - $Y|x \sim \text{Poisson}(\mu)$
 - $E(Y|x) = \mu$
 - $\log(\mu) = a + bx$
 - with $V(Y|x) = \mu$
- Estimating effects
 - $\exp(a) = \mu_0 \Rightarrow \text{Baseline incidence rate}$
 - $\exp(b) = \mu/\mu_0 \Rightarrow \text{Relative risk (RR) for 1 unit increase of } x$
 - $(RR-1) \times 100\% \Rightarrow \text{Percentage risk increase}$

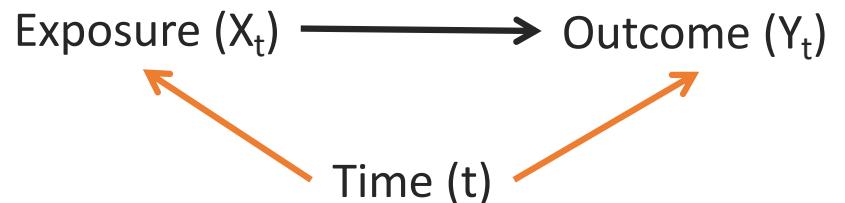
Model assumptions

- Overdispersion
 - In Poisson regression we often find data with $V(E|x)$ larger than $E(Y|x)$, and it underestimates the standard errors
 - Usual solution is to use quasi-Poisson, where $V(Y|x) = \phi\mu$
- Residual autocorrelation
 - The key assumption for regression models is that the outcome measures are independent ($e_t \sim \text{iid}$)
 - In time-series, closer events tend to be more similar than those further apart in time
 - It also affects standard errors and should be accounted to make valid inferences

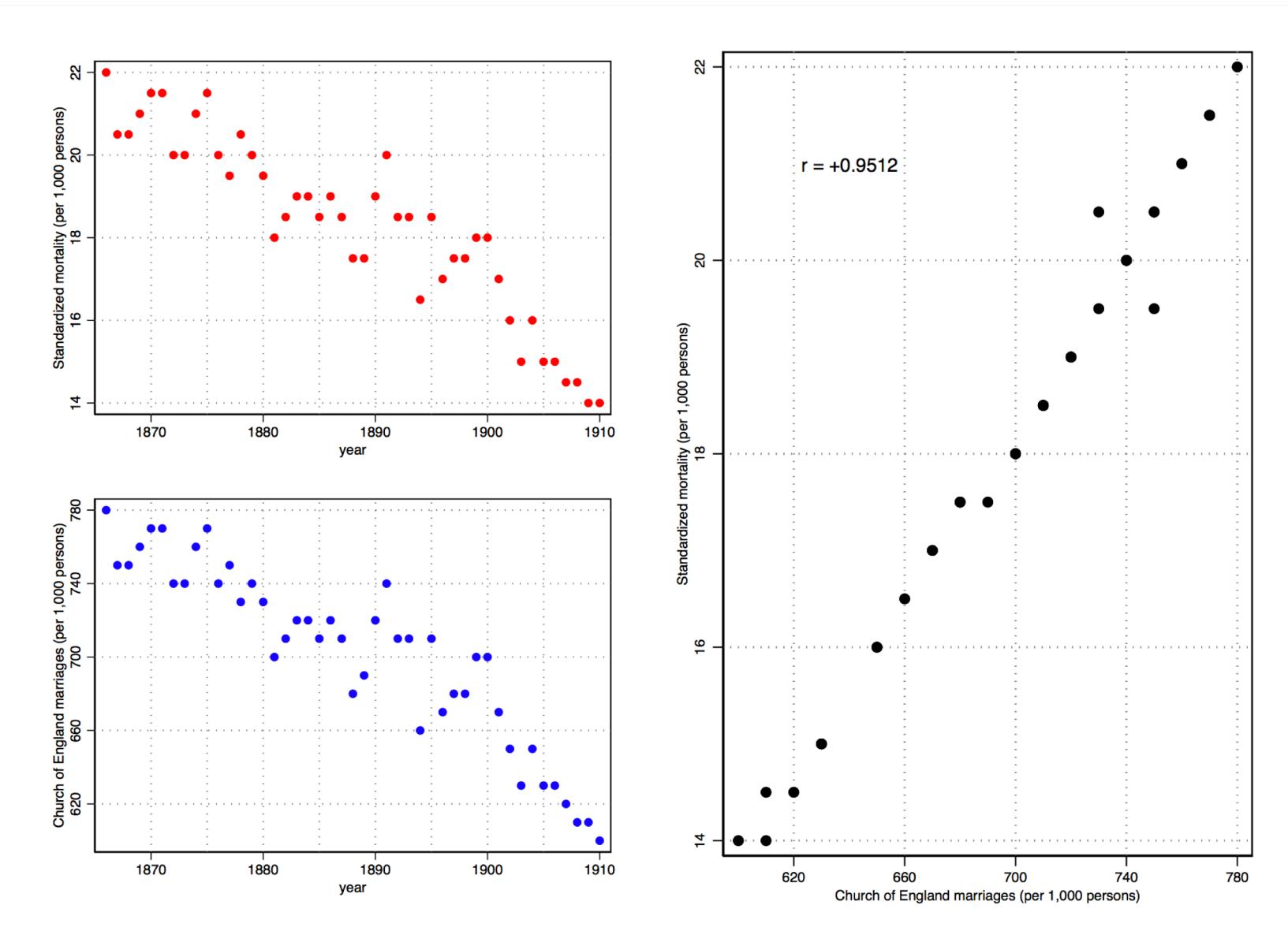


Confounding

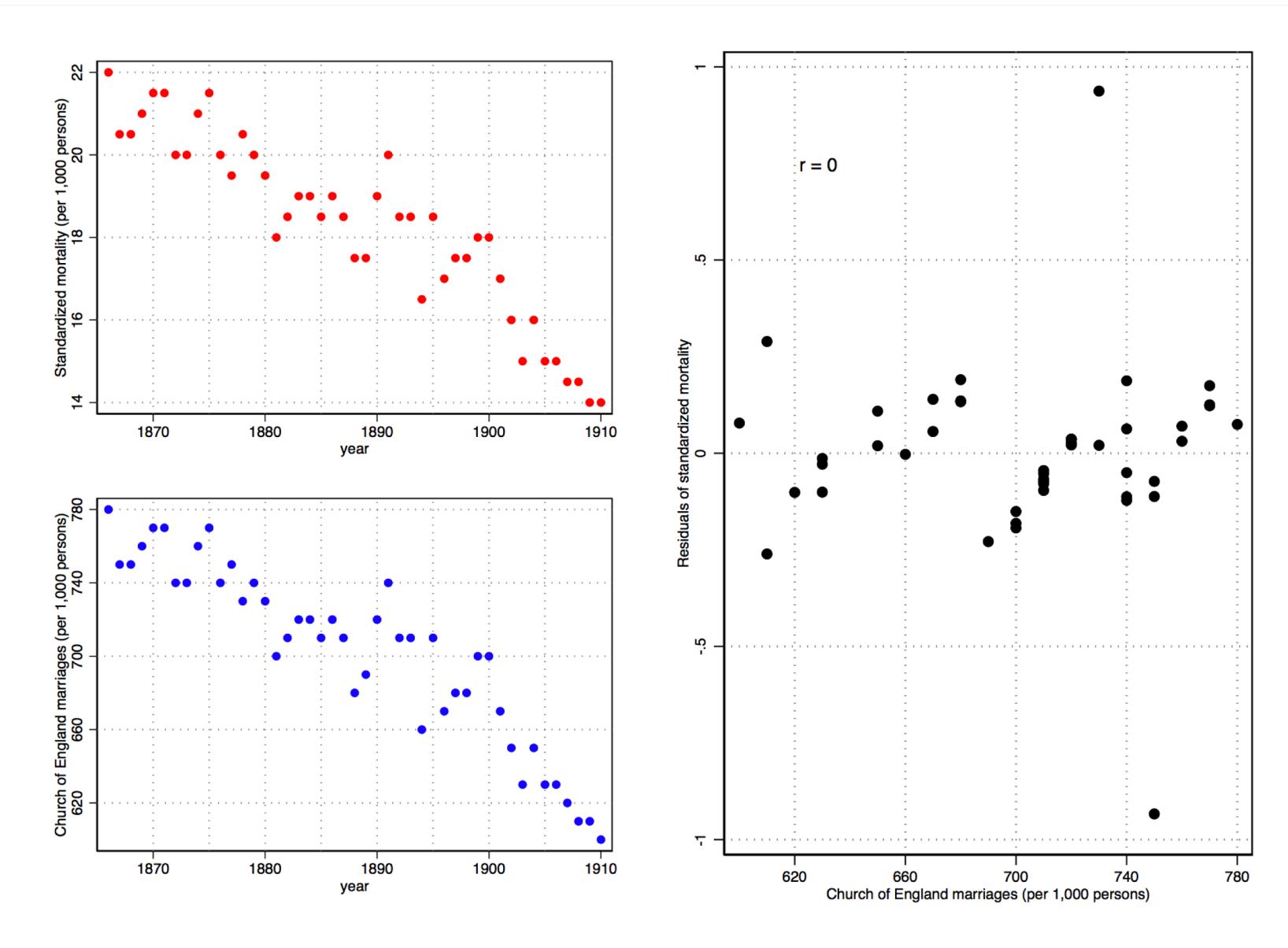
- It must be associated with the exposure (X) being investigated
- It must be independently associated with the outcome (Y) being investigated
- It must not be on the causal pathway between exposure (X) and outcome (Y)



Yule GU. Why do we sometimes get nonsense correlations between time series? J Royal Stat Soc Sci. 1926;89:1-64.



Yule GU. Why do we sometimes get nonsense correlations between time series? J Royal Stat Soc Sci. 1926;89:1-64.

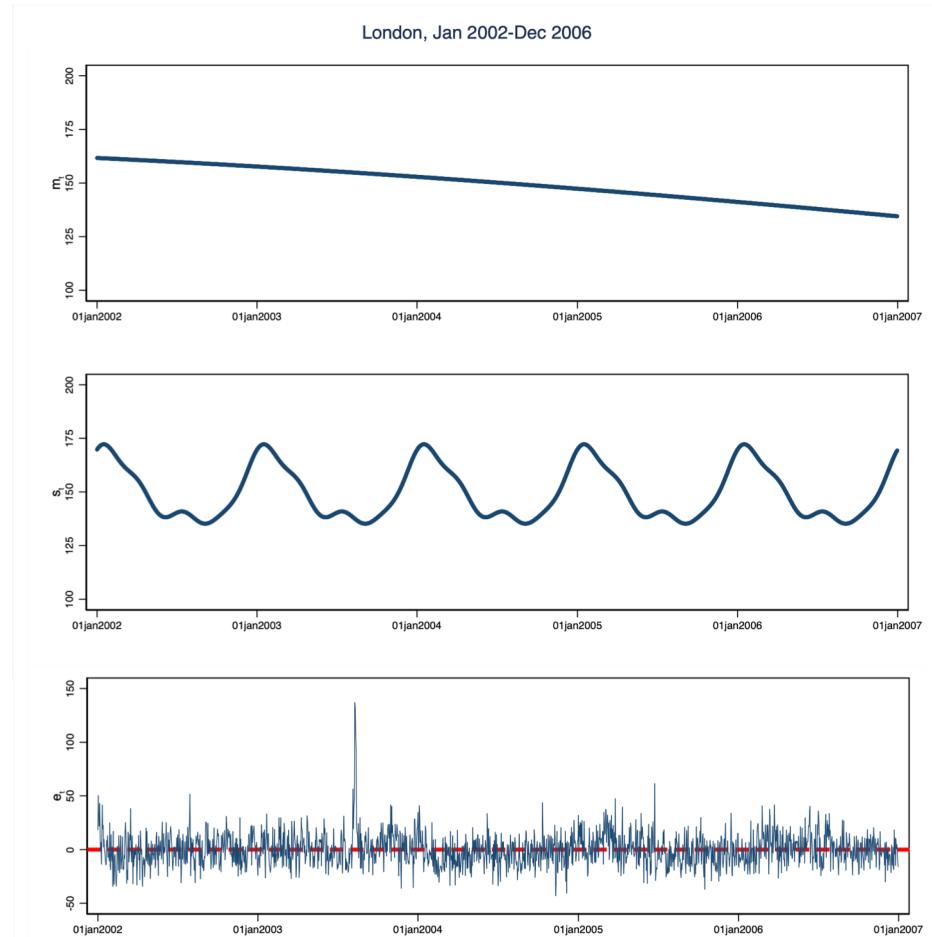


Modelling framework

- Temporal decomposition

$$Y_t = m_t + s_t + e_t$$

- With m_t and s_t as time components (long trend and seasonality) and e_t as residual series
- Underlying trends are **filtered-out** from the time series, allowing the inspection of associations at shorter time scale

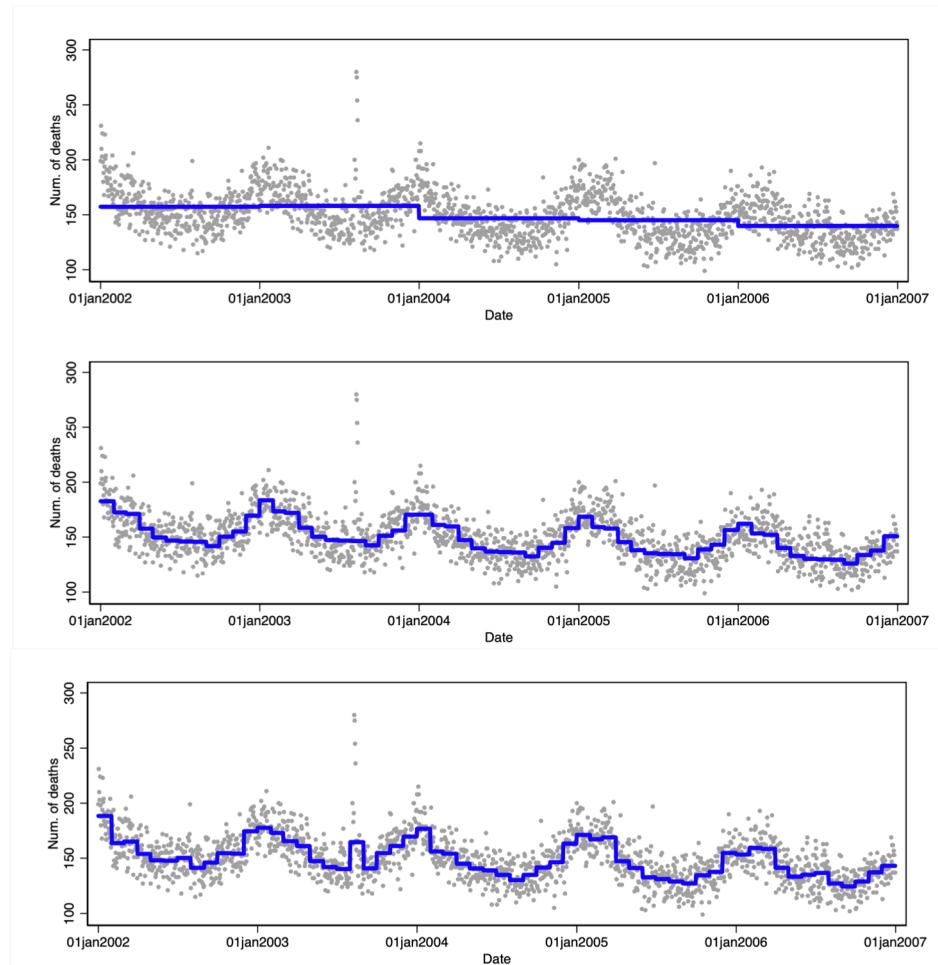


Time-stratified model

- Split the study period into time-intervals estimating a different baseline mortality risk
- Use of indicator variables for year and month

$$Y_t = a + \sum b_i \text{year}_i + \sum d_j \text{month}_j$$

- *Easy to understand, and often captures main long-term patterns*
- *Implicitly assumes biologically implausible jumps in risk between adjacent months*



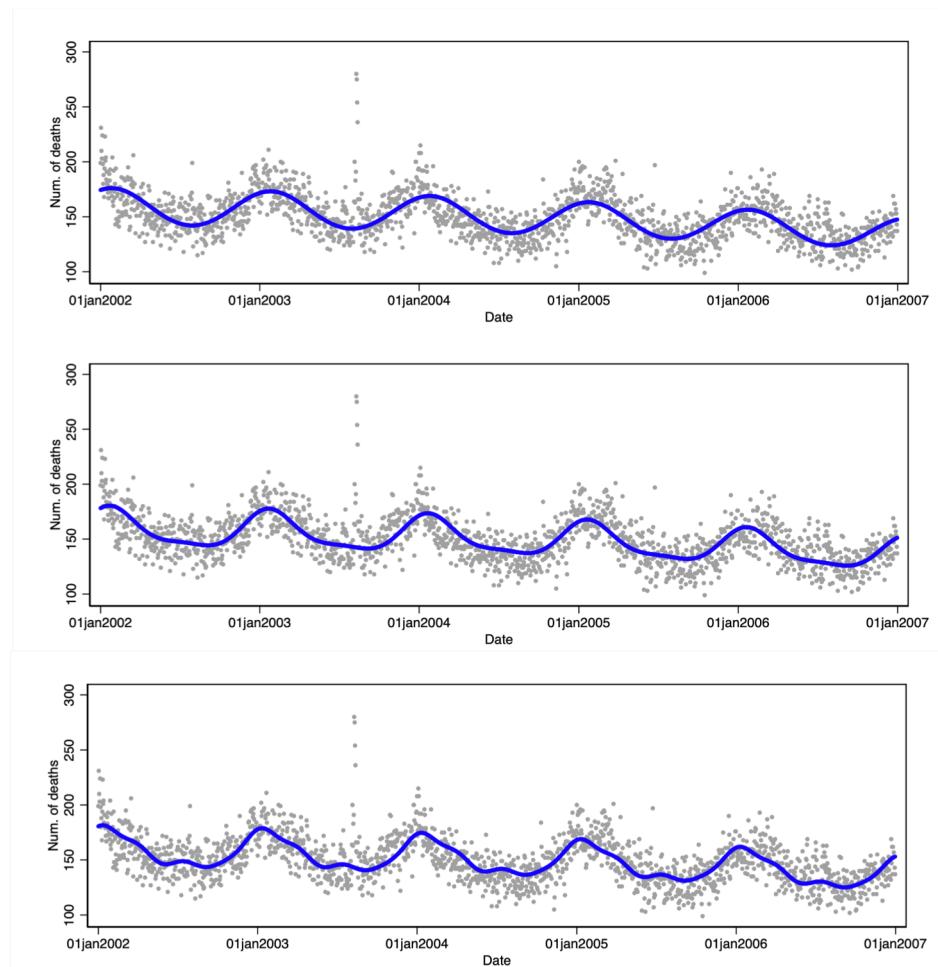
(Bhaskaran *et al.* 2013)

Periodic functions

- Fourier terms (pairs of sine and cosine functions of time) to model seasonal variation in the outcome as a regular wave each year

$$Y_t = a + b_t + \sum d_k \sin(k\pi t / T) + \sum g_k \cos(k\pi t / T)$$

- Suitable to capture very regular seasonal patterns*
- The modelled seasonal pattern is forced to be the same for each year, which may not reflect the data well*



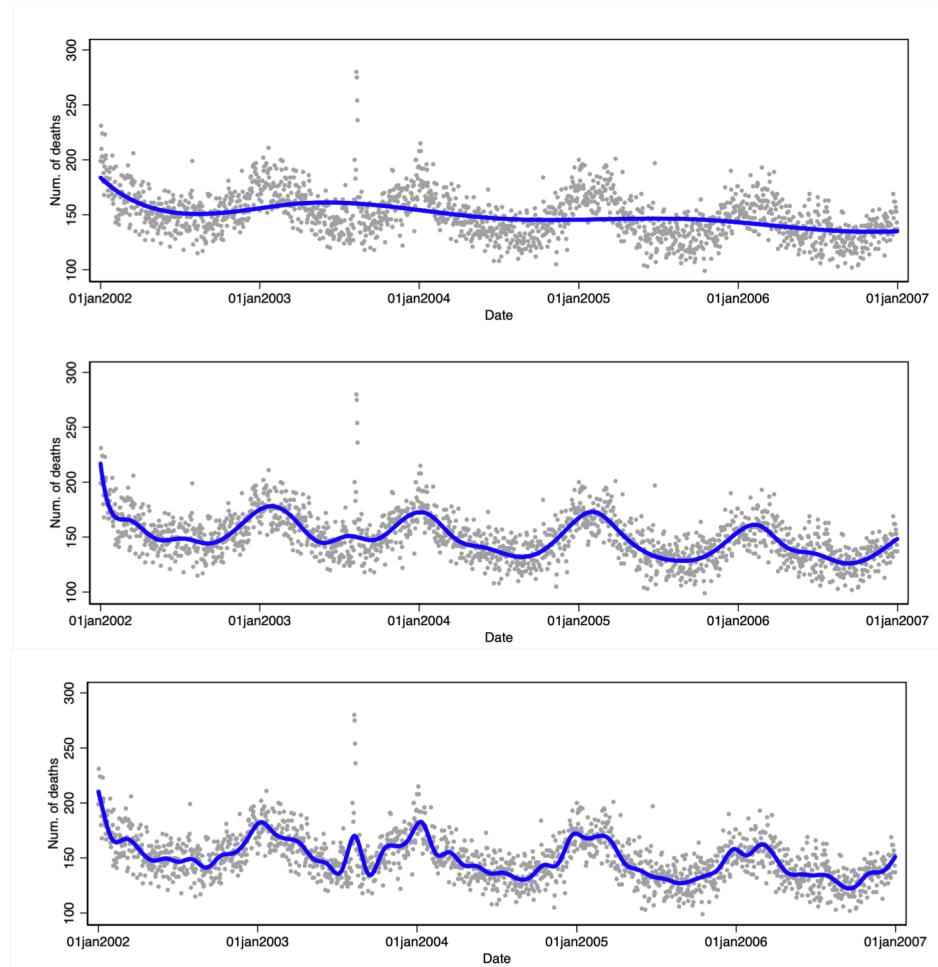
(Bhaskaran *et al.* 2013)

Spline functions

- A set of polynomial curves (commonly cubic) that are joined smoothly end-to-end to cover the study period
- Basis variables are functions of the calendar time

$$Y_t = a + f(t)$$

- Capture seasonal patterns in a way allowed to vary from each year
- It is necessary to decide how many knots there should be, which governs how flexible the curve will be



(Bhaskaran *et al.* 2013)

Comparing modelling strategies

- Time-stratified model

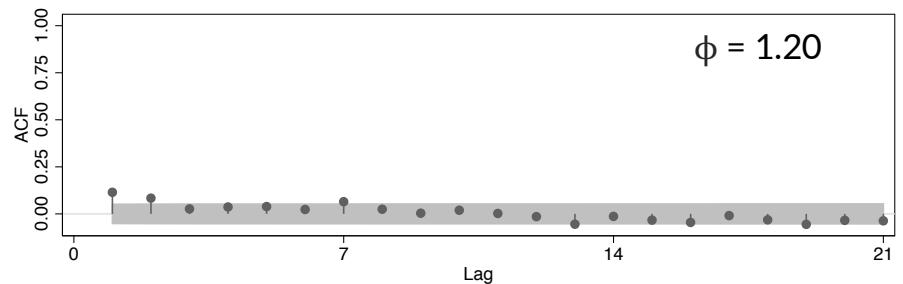
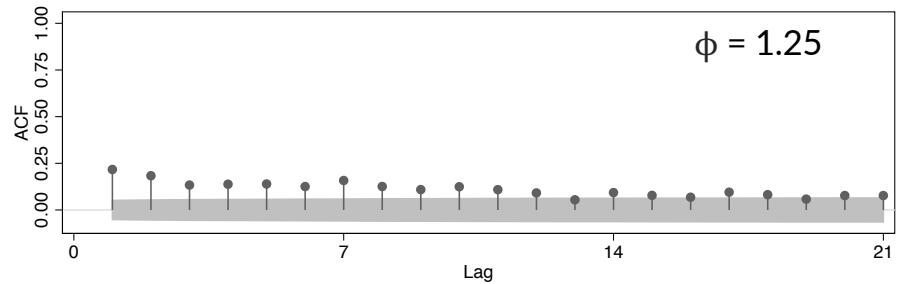
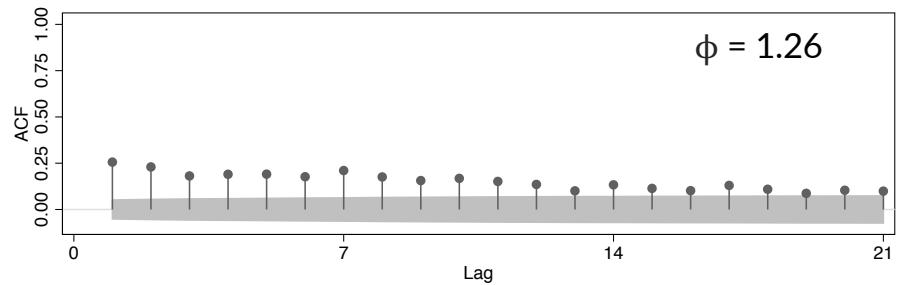
$$Y_t = a + \sum b_i \text{year}_i + \sum d_j \text{month}_j$$

- Periodic functions

$$Y_t = a + bt + \sum d_k \sin(k\pi t/T) + \sum g_k \cos(k\pi t/T)$$

- Spline functions

$$Y_t = a + f(t)$$



Criteria for model selection

- How to select the “right” model among the alternatives? (e.g., smoothing degree for time-trend)
- **Regression diagnostics** (e.g., residuals) might be of help in identifying deviations from model assumptions
- **Statistical criteria**, such as information indices (e.g., log likelihood, AIC), minimizing autocorrelation
- None has been proved as a general reliable method. **A priori assumptions based on epidemiological literature** and sensitivity analysis are recommended

Summary

- Time-series studies provide evidence on short-term associations between environmental exposures and health outcomes
- Time-series regression is similar in principle to any regression analysis but with some specific features
 - Residual autocorrelation
 - Control for time-trends and seasonality

Useful references

- Zeger et al. [On time series analysis of public health and biomedical data.](#) **Annu Rev Public Health** 2006;27:57-79
- Bhaskaran et al. [Time series regression studies in environmental epidemiology.](#) **Int J Epidemiol** 2013;4:1187-95
- Peng and Dominici. [Statistical Methods for Environmental Epidemiology with R – A Case Study in Air Pollution and Health.](#) Springer: New York, 2008.