# Appendix

As the methods used in our analysis are rarely applied in medical research, we provide a short overview of the essentials in this Appendix, a comprehensive account can be found in Helfenstein [1]. We performed a so called intervention analysis, which is a special case of a transfer function analysis. The basic idea beyond a transfer function is to display a time series $y_t$ (number of deaths per day in our case) as the sum of a noise series $n_t$ and a transfer series $u_t$.

$$y_t = n_t + u_t$$

The noise series $n_t$ is the part of the observed time series that cannot be explained by the effect of the input series and captures all of the typical features of a time series such as trends, seasonality, or autocorrelation. In an intervention analysis the transfer series $u_t$ consists of a binary input series $I_t$, which is set to one when a certain event (a soccer match in our case) takes place on day t and zero on every other day, and a so called pulse $w_0$.

$$u_t = w_0 \cdot I_t$$

When there is no match on day t, the time series is, $y_t = n_t + w_0 \cdot 0 = n_t$ while the series can be written as $y_t = n_t + w_0 \cdot 1 = n_t + w_t$ on match days. Thus, the pulse $w_0$ is simply the excess number of deaths on a match day.

The noise series $n_t$ can in most cases (and also in our work) be adequately modelled by a SARIMA-model (seasonal ARIMA model), which is expressed by ARIMA(p,d,q)(P,D,Q)$_s$. To better understand this expression, let us first consider an ARIMA(p,d,q) model and leave the (P,D,Q)$_s$ part for later discussion. Basically, an ARIMA model consists of three terms which model the order of the AR (autoregressive, p), the I (integrated, d), and the MA (moving average, q) part of $n_t$.

**(1) The AR part:**

For convenience, we start with an AR(1) model, that is, an AR model with p=1. This can be interpreted like an ordinary linear regression equation with $y_t$ as the response, $\phi$ as the regression parameter for the single covariate $y_{t-1}$, and a normally distributed random error $e_t$.

$$y_t = \phi y_{t-1} + \epsilon_t$$

The current value of the model, $y_t$ is the sum of the previous value $y_{t-1}$ (multiplied by $\phi$) and the random error. The association between $y_t$ and $y_{t-1}$ is controlled by the AR(1) parameter $\phi$: the larger $\phi$, the higher is the correlation between $y_t$ and $y_{t-1}$. The idea of regressing the current value on its own predecessor explains the term autoregressive for this model. Autoregressive models of higher orders (AR(p) models) are straightforward extensions of the AR(1) process, including the p previous values of the process in the model equation.

$$y_t = \phi_1 y_{t-1} + \phi_2 y_{t-2} + \cdots + \phi_p y_{t-p} + \epsilon_t$$

**(2) The MA part:**

The idea of an MA model is similar to that of an AR model; however, now the current value of the time series, $y_t$, is assumed to depend only on random fluctuations. If random fluctuations on the same day ($\varepsilon_t$), and on the day before ($\varepsilon_{t-1}$) are taken into account, a MA(1) model is defined.

$$y_t = \epsilon_t + \theta \epsilon_{t-1}$$

Higher orders q of an MA model are straightforwardly defined as

$$y_t = \epsilon_t + \theta \epsilon_{t-1} + \ldots + \theta_q \epsilon_{t-q}$$

**(3) The I part**

A time series is said to be stationary, if the mean of the time series does not depend upon time, but is constant throughout the complete time course. For a valid intervention analysis, the noise series $n_t$ has to be stationary. The easiest way to achieve stationarity is by differentiating the time

series by the preceding value s time points ago, where s is the length of the period. Further differentiations with different lags are possible (with d measuring the number of differentiations), but were not necessary in our case.

In the SARIMA model, the seasonal aspect of the noise series $n_t$ is additionally (to the former ARIMA(p,d,q) part) expressed by an ARIMA(P,D,Q)$_s$ term. This seasonal term is assumed to be another ARIMA model with own orders P, D, Q, a seasonal lag parameter s, and own model parameters $\Phi$ and $\Omega$ (now written as capital letters).

Actual model fitting thus involves finding the optimal orders of the SARIMA model, the respective parameters, and the parameter of actual interest, the pulse $w_0$. Box/Jenkins [2] proposed an algorithm for this model identification which is frequently used in applied research. This algorithm consists of four steps:

**(1) Make the original time series ($y_t$) stationary**

Stationarity can be checked via the Dickey-Fuller test. An underlying trend or seasonality is assessable via the empirical autocorrelation and partial autocorrelation functions of $y_t$ at various lags. The autocorrelation at lag k is the correlation of the value $y_t$ and its predecessors $y_{t-k}$. The partial autocorrelation at lag k adjusts for the influence of time points lying between the value $y_t$ and its predecessor $y_{t-k}$, leaving only the adjusted correlation between the two values.

In our case, we found a seasonality of seven days in all 15 models, thus SARIMA models that were differentiated with a lag of seven were fitted (ARIMA(p,0,q)(P,1,Q)$_7$).

**(2) Find a preliminary order of the model**

In the second step, a preliminary order of the ARIMA(p,0,q)(P,1,Q)$_7$ model (that is, p,q, P, and Q) is identified by again referring to the autocorrelation functions.

**(3) Estimate the coefficients of the model**

The coefficients $\theta$, $\varphi$, $\Phi$ and $\Omega$ are estimated by maximum likelihood.

**(4) Check the model by assessing the autocorrelations of the residuals**

As a final step, the adequacy of the model from step (3) has to be evaluated, which is done by demanding no relevant autocorrelations of the residuals. If there are no autocorrelations, the model can be regarded as properly modeling the noise series $n_t$. Finally, all coefficients ($\theta$, $\varphi$, $\Phi$ and $\Omega$) from the noise series are estimated again by maximum likelihood, but now simultaneously with the pulse $w_0$.

**The problem of annual seasonality and model identification**

Referring to the data plot (Figure 1 in the main text), an annual cycle seems to be apparent. Nevertheless, we did not differentiate the observed time series by this period or include an additional seasonal ARIMA term in the model. This was because of the following reasons:

- The Akaike information criterion indicated a worsening of the model, when the time series was differentiated by a period of 365 days and/or when a seasonal term reflecting this periodicity was included in the model.

- Rinne/Specht [3] suggest a maximum lag of $K \approx 2\sqrt{T}$ (in our case the lag length is about 140) when the residuals are checked for autocorrelations as implemented in the method described by Box/Jenkins. Therefore, a lag of 365 days is out of the range to be considered.

- The autocorrelations of the residuals did not noticeably improve after consideration of a seasonal term of 365 days. Thus, we assumed that a SARIMA model without seasonality of 365 days modeled the time series adequately.

**References**

1. Helfenstein U. Box-Jenkins modelling in medical research. Stat Methods Med Res. 1996;5:3-22.

2. Box GEP JGM (1976): Time series analysis: Forecasting and control. San Francisco: Holden-Day [Holden-Day series in time series analysis and digital processing].

3. Rinne H, Specht K Zeitreihen: Statistische Modellierung, Schätzung und Prognose [Time Series: Statistical Modelling, Estimation, Prognosis]. Munich: Vahlen, 2002.