



## Formulating State Space Models in R with Focus on Longitudinal Regression Models

Claus Dethlefsen

Aalborg Hospital, Aarhus University Hospital

Søren Lundbye-Christensen

Aalborg University

---

### Abstract

We provide a language for formulating a range of state space models with response densities within the exponential family. The described methodology is implemented in the R-package `sspir`. A state space model is specified similarly to a generalized linear model in R, and then the time-varying terms are marked in the formula. Special functions for specifying polynomial time trends, harmonic seasonal patterns, unstructured seasonal patterns and time-varying covariates can be used in the formula. The model is fitted to data using iterated extended Kalman filtering, but the formulation of models does not depend on the implemented method of inference. The package is demonstrated on three datasets.

*Keywords:* dynamic models, exponential family, generalized linear models, iterated extended Kalman smoothing, Kalman filtering, seasonality, time series, trend.

---

## 1. Introduction

Generalized linear models, see [McCullagh and Nelder \(1989\)](#), are used when analyzing data where response-densities are assumed to belong to the exponential family. Time series of counts may adequately be described by such models. However, if serial correlation is present or if the observations are overdispersed, these models may not be adequate, and several approaches can be taken. The book by [Diggle, Heagerty, Liang, and Zeger \(2002\)](#) gives an excellent review of many approaches incorporating serial correlation and overdispersion in generalized linear models. Dynamic generalized linear models (DGLM), often called state space models, also address those problems and are treated in a paper by [West, Harrison, and Migon \(1985\)](#) in a conjugate Bayesian setting. They have been subject to further research by e.g. [Zeger \(1988\)](#) using generalized estimating equations (GEE), [Gamerman \(1998\)](#) using Markov chain Monte Carlo (MCMC) methods and [Durbin and Koopman \(1997\)](#) using iterated extended Kalman filtering and importance sampling.

Standard statistical software does not include procedures for DGLMs and only sparse support for Gaussian state space models. There is a need for a simple, yet flexible way of specifying complicated non-Gaussian state space models. Often, one needs to tailor make software for each specific application in mind. A function, `StructTS`, has been developed for analysis of a subclass of Gaussian state space models, see Ripley (2002). The binary library `SsfPack` for Ox may be used freely for academic research and provides a tool set for analysis of Gaussian state space models with some support for non-Gaussian models, see Koopman, Shephard, and Doornik (1999). The interface is very flexible, but not as easy to use as a `glm` call in R.

Section 2 describes Gaussian state space models and shows how generalized linear models can naturally be extended to allow the parameters to evolve over time. We define components (e.g. trend and seasonal components) that separate the time series into parts that may be inspected individually after analysis. In Section 3 the syntax for defining objects describing the proposed state space models are described as a simple, yet powerful, extension to the `glm`-call in R (R Development Core Team 2006). The techniques are illustrated on three examples in Section 4.

## 2. State space models

The Gaussian state space model for univariate observations involves two processes, namely the state process (or latent process),  $\{\boldsymbol{\theta}_k\}$ , and the observation process,  $\{y_k\}$ . The random variation in the state space model is specified through descriptions of the sampling distribution, the evolution of the state vector, and the initialization of the state vector.

Let  $\{y_k\}$  be measured at timepoints  $t_k$  for  $k = 1, \dots, n$ . The state space model is defined by

$$y_k = \mathbf{F}_k^\top \boldsymbol{\theta}_k + \nu_k, \quad \nu_k \sim \mathcal{N}(0, V_k) \quad (1)$$

$$\boldsymbol{\theta}_k = \mathbf{G}_k \boldsymbol{\theta}_{k-1} + \boldsymbol{\omega}_k, \quad \boldsymbol{\omega}_k \sim \mathcal{N}_p(\mathbf{0}, \mathbf{W}_k) \quad (2)$$

$$\boldsymbol{\theta}_0 \sim \mathcal{N}_p(\mathbf{m}_0, \mathbf{C}_0). \quad (3)$$

We assume that the disturbances  $\{\nu_k\}$  and  $\{\boldsymbol{\omega}_k\}$  are both serially independent and also independent of each other. The possible time-dependent quantities  $\mathbf{F}_k$ ,  $\mathbf{G}_k$ ,  $V_k$  and  $\mathbf{W}_k$  may depend on a parameter vector, but this is suppressed in the notation.

We now consider the case where the state process is Gaussian and the sampling distribution belongs to the exponential family,

$$p(y_k | \eta_k) = \exp \{y_k \eta_k - b_k(\eta_k) + c_k(y_k)\}. \quad (4)$$

The density (4) contains the Gaussian, Poisson, gamma and the binomial distributions as special cases. The natural parameter  $\eta_k$  is related to the linear predictor  $\lambda_k$  by the equation  $\eta_k = v(\lambda_k)$  or equivalently  $\lambda_k = u(\eta_k)$ . The linear predictor in a generalized linear model is of the form  $\lambda_k = \mathbf{Z}_k \boldsymbol{\beta}$ , where  $\mathbf{Z}_k$  is a row vector of explanatory variables and  $\boldsymbol{\beta}$  is the vector of regression parameters. The link function,  $g$ , relates the mean,  $\mathbb{E}(y_k) = \mu_k$ , and the linear predictor,  $\lambda_k$ , as  $g(\mu_k) = \lambda_k$ . The inverse link function,  $h$ , is defined as  $\mu_k = \tau(\eta_k) = h(\lambda_k)$ , where  $\tau$  is the mean value mapping. The following relations hold  $\eta_k = v(\lambda_k) = \tau^{-1}(h(\lambda_k))$  and  $\lambda_k = u(\eta_k) = g(\tau(\eta_k))$ , where  $u$  is the inverse of  $v$ . The link function is said to be canonical if  $\eta_k = \lambda_k$ , i.e. if  $g = \tau^{-1}$ .

## 2.1. Dynamic extension

The static generalized linear model is extended by adding a dynamic term,  $\mathbf{X}_k\boldsymbol{\beta}_k$ , to the linear predictor, where  $\boldsymbol{\beta}_k$  is varying randomly over time according to a first order Markov process. Hence,

$$\lambda_k = \mathbf{Z}_k\boldsymbol{\beta} + \mathbf{X}_k\boldsymbol{\beta}_k, \quad (5)$$

where  $\boldsymbol{\beta}$  is the coefficient of the static component and  $\{\boldsymbol{\beta}_k\}$  are the time-varying coefficients of the dynamic component.

For notational convenience, we will use the notation

$$\lambda_k = \mathbf{F}_k^\top \boldsymbol{\theta}_k, \quad \boldsymbol{\theta}_k = \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\beta}_k \end{pmatrix}. \quad (6)$$

The evolution through time of the state vector,  $\boldsymbol{\theta}_k$ , is modelled by the relation

$$\boldsymbol{\theta}_k = \mathbf{G}_k\boldsymbol{\theta}_{k-1} + \boldsymbol{\omega}_k, \quad (7)$$

for an evolution matrix  $\mathbf{G}_k$ , determined by the model. The error terms,  $\{\boldsymbol{\omega}_k\}$ , are assumed to be independent Gaussian variables with zero mean and variance  $\text{VAR}(\boldsymbol{\omega}_k)$ , with non-zero entries corresponding to the entries of the time-varying coefficients,  $\boldsymbol{\beta}_k$ , and zero elsewhere.

The model is fully specified by the initializing parameters  $\mathbf{m}_0$  and  $\mathbf{C}_0$ , the matrices  $\mathbf{F}_k$ ,  $\mathbf{G}_k$ , and the variance parameters  $V_k$  and  $\text{VAR}(\boldsymbol{\omega}_k)$ . The variances may be parametrized as e.g.  $\text{VAR}(\boldsymbol{\omega}_k) = \psi \cdot \text{diag}(1, 0, 0, 1, 1)$  or  $\text{VAR}(\boldsymbol{\omega}_k) = \text{diag}(\psi_1, \psi_2, \psi_2)$ .

## 2.2. Inferential procedures

For a Gaussian state space model, we write  $\boldsymbol{\theta}_k|\mathcal{D}_k \sim \mathcal{N}_p(\mathbf{m}_k, \mathbf{C}_k)$ , where  $\mathcal{D}_k$  is all information available at time  $t_k$ . The Kalman filter recursively yields  $\mathbf{m}_k$  and  $\mathbf{C}_k$  with the recursion starting in  $\boldsymbol{\theta}_0 \sim \mathcal{N}_p(\mathbf{m}_0, \mathbf{C}_0)$ .

Assessment of the state vector,  $\boldsymbol{\theta}_k$ , using all available information,  $\mathcal{D}_n$ , is called Kalman smoothing and we write  $\boldsymbol{\theta}_n|\mathcal{D}_n \sim \mathcal{N}_p(\tilde{\mathbf{m}}_n, \tilde{\mathbf{C}}_n)$ . Starting with  $\tilde{\mathbf{m}}_n = \mathbf{m}_n$  and  $\tilde{\mathbf{C}}_n = \mathbf{C}_n$ , the Kalman smoother is a backwards recursion in time,  $k = n - 1, \dots, 1$ .

For exponential family sampling distributions, *the iterated extended Kalman filter* yields an approximation to the conditional distribution of the state vector given  $\mathcal{D}_n$ , see e.g. [Durbin and Koopman \(2000\)](#). By Taylor expansion, the sample distribution (4) is approximated with a Gaussian density, giving an approximating Gaussian state space model. The conditional distribution of the state vector given  $\mathcal{D}_n$  in the exact model and in the Gaussian approximation have the same mode. The iterated extended Kalman filter is used as filter and smoother method in `sspir`.

## 2.3. Decomposition

The variation in the linear predictor, random or not, may be decomposed into four components: a time trend ( $T_k$ ), harmonic seasonal patterns ( $H_k$ ), unstructured seasonal patterns ( $S_k$ ), and a regression with possibly time-varying covariates ( $R_k$ ).

Each component may contain static and/or dynamic components, which is specified by zero and non-zero diagonal elements in  $\text{VAR}(\boldsymbol{\omega}_k)$ , respectively, as described in the following.

The block-diagonal evolution matrix takes the form

$$\mathbf{G}_k = \begin{bmatrix} \mathbf{G}_k^{(1)} & & & \\ & \mathbf{I} & & \\ & & \mathbf{G}_k^{(3)} & \\ & & & \mathbf{I} \end{bmatrix},$$

where  $\mathbf{G}_k^{(1)}$  is defined in (9), and  $\mathbf{G}_k^{(3)}$  in (12). The components are only present if the model includes the corresponding terms.

The linear predictor,

$$\begin{aligned} \lambda_k &= \mathbf{T}_k \boldsymbol{\theta}_k^{(1)} + \mathbf{H}_k \boldsymbol{\theta}_k^{(2)} + \mathbf{S}_k \boldsymbol{\theta}_k^{(3)} + \mathbf{R}_k \boldsymbol{\theta}_k^{(4)} \\ &= T_k + H_k + S_k + R_k. \end{aligned}$$

will be detailed in the following.

### *Time trend*

The long term trend is usually modelled by a sufficiently smooth function. In static regression models, this can be done by e.g. a high degree polynomial, a spline, or a generalized additive model. In the dynamic setting, however, a low degree polynomial with time-varying coefficient may suffice.

By stacking a polynomial,  $q(t) = b_0 + b_1 t + \dots + b_p t^p$ , and the first  $p$  derivatives, the transition from  $t_{k-1}$  to  $t_k$  obeys the relation

$$\begin{bmatrix} q(t_k) \\ q'(t_k) \\ \vdots \\ q^{(p)}(t_k) \end{bmatrix} = \mathbf{G}_k^{(1)} \begin{bmatrix} q(t_{k-1}) \\ q'(t_{k-1}) \\ \vdots \\ q^{(p)}(t_{k-1}) \end{bmatrix}, \quad (8)$$

where  $\Delta t_k = t_k - t_{k-1}$ , and the upper triangular transition matrix is given by

$$\mathbf{G}_k^{(1)} = \begin{bmatrix} 1 & \Delta t_k & \dots & \Delta t_k^p / p! \\ & 1 & \dots & \Delta t_k^{p-1} / (p-1)! \\ & & \ddots & \vdots \\ & & & 1 \end{bmatrix}. \quad (9)$$

Using  $\boldsymbol{\theta}_k^{(1)}$  for the left hand side of (8), a polynomial growth model with time-varying coefficients can be written as  $\boldsymbol{\theta}_k^{(1)} = \mathbf{G}_k^{(1)} \boldsymbol{\theta}_{k-1}^{(1)} + \boldsymbol{\omega}_k^{(1)}$ . The error term has variance  $\text{VAR}(\boldsymbol{\omega}_k^{(1)}) = \Delta t_k \mathbf{W}^{(1)}$ , where  $\mathbf{W}^{(1)}$  is diagonal in the case with independent random perturbations in each of the derivatives.

The trend component is the first element in  $\boldsymbol{\theta}_k^{(1)}$ , i.e.

$$T_k = \mathbf{T}_k \boldsymbol{\theta}_k^{(1)} = [1 \ 0 \ \dots \ 0] \boldsymbol{\theta}_k^{(1)}.$$

Alternatively, the time trend may be modelled as a random function,  $q(t)$ , for which the increments over time are described by a random walk, resulting in a cubic spline, see [Kitagawa](#)

and Gersch (1984). The transition is the same as in (8) with  $p = 2$ , but only one variance parameter is necessary as,

$$\text{VAR}(\boldsymbol{\omega}_k^{(1)}) = \sigma_w^2 \begin{bmatrix} \Delta t_k^3/3 & \Delta t_k^2/2 \\ \Delta t_k^2/2 & \Delta t_k \end{bmatrix}. \quad (10)$$

### Harmonic seasonal pattern

Seasonal patterns with a given period,  $m$ , can be described by the following  $d$ th degree trigonometric polynomial

$$\begin{aligned} H_k &= \mathbf{H}_k \boldsymbol{\theta}_k^{(2)} \\ &= \sum_{i=1}^d \left\{ \theta_{c,i} \cos \left( i \cdot \frac{2\pi}{m} t_k \right) + \theta_{s,i} \sin \left( i \cdot \frac{2\pi}{m} t_k \right) \right\} \\ &= [c_{1k} \ \cdots \ c_{dk} \ s_{1k} \ \cdots \ s_{dk}] \boldsymbol{\theta}_k^{(2)}, \end{aligned} \quad (11)$$

where  $c_{ik} = \cos(i \cdot 2\pi t_k / m)$  and  $s_{ik} = \sin(i \cdot 2\pi t_k / m)$ . This component can be used to describe seasonal effects showing cyclic patterns. Further seasonal components may be added for each period of interest.

The random fluctuations in  $\boldsymbol{\theta}_k^{(2)}$  is modelled by a random walk,  $\boldsymbol{\theta}_k^{(2)} = \boldsymbol{\theta}_{k-1}^{(2)} + \boldsymbol{\omega}_k^{(2)}$  with  $\text{VAR}(\boldsymbol{\omega}_k^{(2)}) = \Delta t_k \mathbf{W}^{(2)}$ .

### Unstructured seasonal component

For equidistant observations, a commonly used parameterization for the seasonal component is to let the effects,  $\gamma_k$ , for each period sum to zero in the static case, or to a white noise error sequence in the time-varying case, see Kitagawa and Gersch (1984). For an integer period,  $m$ , the sum-to-zero constraint can be expressed as  $\sum_{i=0}^{m-1} \gamma_{k-i} = 0$  in the static case, and in the dynamic case,  $\sum_{i=0}^{m-1} \gamma_{k-i} = \omega_k^{(3)}$ , with  $\omega_k^{(3)} \sim \mathcal{N}(0, \sigma_w^2)$ . This is expressed in matrix form by letting  $\boldsymbol{\theta}_k^{(3)} = [\gamma_k, \gamma_{k-1}, \dots, \gamma_{k-m+2}]^\top$ , and defining the  $(m-1) \times (m-1)$  matrix

$$\mathbf{G}_k^{(3)} = \begin{bmatrix} -1 & -1 & \cdots & -1 \\ 1 & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 1 & 0 \end{bmatrix}. \quad (12)$$

Then,  $\boldsymbol{\theta}_k^{(3)} = \mathbf{G}_k^{(3)} \boldsymbol{\theta}_{k-1}^{(3)} + \boldsymbol{\omega}_k^{(3)}$ , with  $\text{VAR}(\boldsymbol{\omega}_k^{(3)}) = \mathbf{W}^{(3)} = \text{diag}(\sigma_w^2, 0, \dots, 0)$  defines the evolution of the seasonal component. The corresponding term in the linear predictor is extracted by

$$S_k = \mathbf{S}_k \boldsymbol{\theta}_k^{(3)} = [1 \ 0 \ \cdots \ 0] \boldsymbol{\theta}_k^{(3)}.$$

### Regression component

Observed time-varying covariates,  $\mathbf{R}_k$ , enter the model through the usual regression term

$$R_k = \mathbf{R}_k \boldsymbol{\theta}_k^{(4)},$$

with  $\boldsymbol{\theta}_k^{(4)} = \boldsymbol{\theta}_{k-1}^{(4)} + \boldsymbol{\omega}_k^{(4)}$  and  $\text{VAR}(\boldsymbol{\omega}_k^{(4)}) = \Delta t_k \mathbf{W}^{(4)}$ . The structure of  $\mathbf{W}^{(4)}$  is specified by the modeller and depends on the context.

### 3. Specification of state space objects

The package `sspir` can be downloaded and installed from <http://CRAN.R-project.org/> and is then activated in R by `library("sspir")`. Assuming that the data are available either in a dataframe or in the current environment, then a state space model is setup using `glm`-style formula and family arguments. Terms are considered static unless embraced by the special function `tvar()`, described further in Section 3.2.

#### 3.1. State space model objects

In `sspir`, a state space model is defined as an object from the class `ssm`. The object defines the model and contains the slots that are needed for the subsequent statistical analysis.

The definition of a state space model object has the following syntax

```
ssm(formula, family=gaussian, data, subset, fit=TRUE,
    phi, m0, C0, Fmat, Gmat, Vmat, Wmat)
```

The call is designed to be similar to the `glm` call. The elements in the call are

**formula** a specification of the linear predictor (5) of the model. The syntax is defined in Section 3.2.

**family** a specification of the observation error distribution and link function to be used in the model, as in a `glm`-call. This can be a character string naming a family function, a family function or the result of a call to a family function. Currently, only Poisson with log-link, binomial with logit-link, and Gaussian with identity-link have been implemented. It is possible to expand with further combinations within the exponential family.

**data** an optional data frame containing the variables in the model. By default the variables are taken from `'environment(formula)'`, typically the environment from which `'ssm'` is called. The response has to be of class `ts`.

**subset** an optional vector specifying a subset of observations to be used in the fitting process.

**fit** a logical defaulting to `TRUE` which means that the iterated extended Kalman smoother is used to fit the model. If `FALSE`, the model is only *defined* and no inferential calculations are made.

**phi** a vector of hyper parameters that are passed directly to `Fmat`, `Gmat`, `Vmat`, and `Wmat`. If `phi` is not provided, it is default set to a vector of ones with the length determined by the number of hyper parameters needed on the basis of the `formula` provided.

**m0** a vector with the initial state vector. Defaults to a vector of zeros.

**C0** a matrix with the variance matrix of the initial state. Defaults to a diagonal matrix with diagonal entries set to  $10^6$ .

**Fmat** a function giving the regression matrix at a given timepoint. If not supplied, this is constructed from the `formula`.

**Gmat** a function giving the evolution matrix at a given timepoint. If not supplied, this is constructed from the `formula`.

**Wmat** a function giving the evolution variance matrix at a given timepoint. If not supplied, this is constructed from the `formula`.

**Vmat** a function giving the observation variance matrix at a given timepoint. If not supplied, this is constructed from the `formula`.

The call creates an object defining the system matrices  $\mathbf{F}_t$ ,  $\mathbf{G}_t$ ,  $\mathbf{W}_t$ , and  $\mathbf{V}_t$  in terms of functions, returning the matrix in question at a given time point. For example, the `Wmat` function could be defined as

```
Wmat <- function(tt,x,phi) {
  if (tt==10) return(matrix(phi[2]))
  else return(matrix(phi[3]))
}
```

Here, `Wmat` has value `phi[3]` at all time-points except time-point `tt==10`, where the value is `phi[2]`. This provides a mechanism of incorporating interventions and change-points at any given time. Note, that the call to `ssm` creates the functions which can be re-used. In the following example, the `Wmat` function is first created in the call to `ssm`, and then manually changed so that the second diagonal entry is larger at timepoint `tt==10`. Finally, the model is fitted using the function `kfs`.

```
gasmodel <- ssm(log10(UKgas) ~ -1 + tvar(polytime(time,1)),fit=FALSE)
Wold <- Wmat(gasmodel)
Wmat(gasmodel) <- function(tt,x,phi) {
  W <- Wold(tt,x,phi)
  if (tt==10) {W[2,2] <- 100*W[2,2]; return(W)}
  else return(W)
}
gasmodel.fit <- kfs(gasmodel)
```

### 3.2. Model formulas

A model formula is built up as in a `glm` call in R. The response appears on the left hand side of a tilde (`~`) and on the right hand side the explanatory variables, factors and continuous variables, appear. However, to specify time-varying regression coefficients, we have defined a special notation, `tvar()`, in which these are enclosed.

For example, the formula

```
y ~ z + tvar(x)
```

will correspond to covariates,  $\mathbf{z}$  and  $\mathbf{x}$ , of which  $\mathbf{z}$  has a static parameter and  $\mathbf{x}$  has a dynamic parameter. An implicit intercept is also included in the model, unless the term `-1` appears in the formula. When `tvar` enters a formula and `-1` is *not* included, the intercept will always be time-varying, *i.e.* a random walk is added to the linear predictor. Thus, this model corresponds to the state space model with the linear predictor specified as  $\lambda_k = \mathbf{Z}_k\boldsymbol{\beta} + \mathbf{X}_k\boldsymbol{\beta}_k$ ,  $\mathbf{Z}_k$  being the  $k$ th row in the  $n \times 1$  matrix  $\mathbf{Z} = [\mathbf{z}]$  and  $\mathbf{X}_k$  the  $k$ th row of the  $n \times 2$  matrix  $\mathbf{X} = [1 \ \mathbf{x}]$ . The R command `model.matrix` applied to the formula `y ~ z + x` yields the  $n \times 3$  matrix  $[1 \ z \ x]$ , in which the rows are  $\mathbf{F}_k^\top$ .

The polynomial time trend, (9), is specified using the function,

```
polytime(time,degree=p)
```

Note that `polytime` is different than the built-in R-function `poly` since the latter produces a design matrix with orthonormal columns.

The harmonic seasonal pattern, (11), is specified using the function,

```
polytrig(time,period=m,degree=d)
```

whereas the unstructured seasonal pattern, (12), is specified using the function,

```
sumseason(time,period=m)
```

Regression components are specified using the usual Wilkinson-Rogers formula notation in R. The model matrix does not contain information about which variables are time-varying. This distinction is implemented by specifying the variance matrix,  $\text{VAR}(\boldsymbol{\omega}_k)$ , with zeros in entries corresponding to static parameters and non-zero entries otherwise.

### 3.3. Inference

When a model has been defined using `ssm` and the option `fit` has been set to `TRUE`, the *iterated extended Kalman smoother* has been applied. The output is stored with the `ssm` object and can be extracted by the function `getFit`. This is a list that contains the estimated mean (as the component `m`), and variance matrices (as the component `C`) of the state vector,  $\boldsymbol{\theta}_k$ , as well as the approximate log-likelihood (as the component `loglik`) based on the Gaussian approximation to the state space model.

If `ssm` has been called with the option `fit` set to `FALSE`, the function `kfs` returns the output object.

The following example defines a state space model, runs the iterated extended Kalman smoother, and finally extracts the fitted information into the variable `vd.fit`. This variable is a list where, `vd.fit$m[t,]` contains the conditional mean  $E[\boldsymbol{\theta}_t|\mathcal{D}_n]$  and `vd.fit$C[[t]]` contains the corresponding variance matrix,  $\text{VAR}[\boldsymbol{\theta}_t|\mathcal{D}_n]$ .

```
vd <- ssm( y ~ tvar(1) + seatbelt + sumseason(time,12),
          family=poisson(link="log"),
          data=vanddrivers,
          phi = c(1,0.0004945505),
          C0=diag(13)*1000
        )
vd.fit <- getFit(vd)
```

## 4. Examples

In this section, three examples of specification and application of state space models will be presented. The examples include Gaussian and Poisson observation densities. The time series are decomposed into components of trend and seasonality and also inclusion of external covariates is illustrated. The main focus will be on formulation of the state space object, how a relevant data analysis can be performed, and how to present the output from the analysis, based on this object.

### Example 4.1 (Gas consumption)

A dataset provided with *R* is the quarterly UK gas consumption from 1960 to 1986, in millions of therms (*Durbin and Koopman 2001*, p. 233). As response, we use the (base 10) logarithm

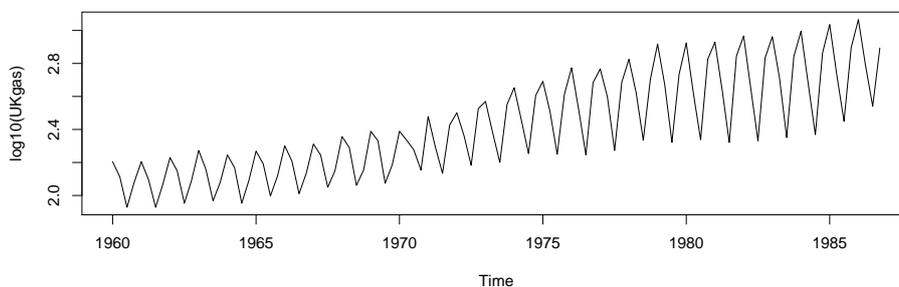


Figure 1: Log-transformed UK gas consumption, recorded quarterly from 1960 to 1986.

of the UK gas consumption (displayed in Figure 1), which we assume is normal distributed. We fit a model with a first order polynomial trend with time-varying coefficients and an unstructured seasonal component, also varying over time.

$$y_k = \log_{10} UKgas_k = T_k + S_k + \varepsilon_k,$$

with the linear trend component being

$$T_k = T_{k-1} + \beta_{k-1} + \omega_k^{(1)}, \quad \beta_k = \beta_{k-1} + \omega_k^{(2)}$$

and the seasonal component being

$$S_k = \gamma_k = -(\gamma_{k-1} + \gamma_{k-2} + \gamma_{k-3}) + \omega_k^{(3)},$$

where  $\omega_k^{(1)}$ ,  $\omega_k^{(2)}$  and  $\omega_k^{(3)}$  are independent noise components. The corresponding vector notation is

$$y_k = [1 \ 0 \ 1 \ 0 \ 0 \ 0] \boldsymbol{\theta}_k + \varepsilon_k,$$

by blocking the evolution matrices of (9) and (12) we get

$$\boldsymbol{\theta}_k = \begin{bmatrix} T_k \\ \beta_k \\ \gamma_k \\ \gamma_{k-1} \\ \gamma_{k-2} \end{bmatrix} = \begin{bmatrix} 1 & 1 & & & \\ 0 & 1 & & & \\ & & -1 & -1 & -1 \\ & & 1 & 0 & 0 \\ & & 0 & 1 & 0 \end{bmatrix} \boldsymbol{\theta}_{k-1} + \begin{bmatrix} \omega_k^{(1)} \\ \omega_k^{(2)} \\ \omega_k^{(3)} \\ 0 \\ 0 \end{bmatrix}.$$

The model is specified, fitted, and plotted in **sspir** by

```
phistart <- c(3.7e-4,0,1.7e-5,7.1e-4)
gasmodel <- ssm(log10(UKgas) ~ -1 + tvar(polytime(time,1)) +
               tvar(sumseason(time,4)), phi=phistart)
fit <- getFit(gasmodel)

plot(fit$m[,1:3])
```

Here, the estimated variances are taken from an external maximum likelihood algorithm provided by the function **StructTS**, Ripley (2002), which is standard in R. The decomposition in trend, slope and season components is displayed in Figure 2. In 1971, the slope increases

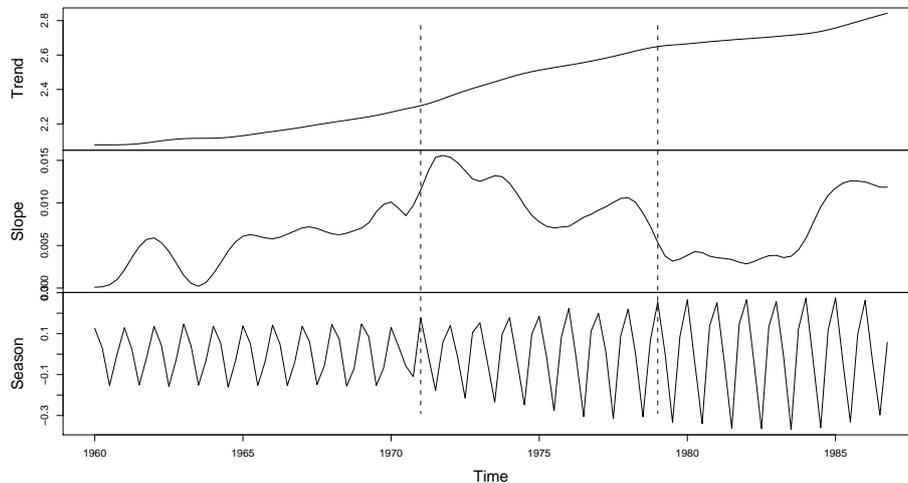


Figure 2: Time-varying trend, slope, and seasonal components in UK gas consumption.

from approximately 0.005 to approximately 0.014 and returns to this level in 1979. At the end of the observation period, the slope increases again. Similarly, it is seen, that the amplitude of the seasonal component is fairly constant from 1960-1971, after which it increases in the period 1971-1979 and then it stabilizes. The analysis can be reproduced in **sspir** by `demo("gas")`.

#### Example 4.2 (Vandriivers)

Let  $y_t$  be the monthly numbers of light goods van drivers killed in road accidents, from January 1969 to December 1984 (192 observations). On January 31st, 1983, a seat belt law was introduced. The interest is to quantify the effect of the seat belt legislation law. For further information about the data set consult Harvey and Durbin (1986).

Here we use a state space model for Poisson data with a 13-dimensional latent process, consisting of an intervention parameter, `seatbelt`, changing value from zero to one in February 1983, a constant monthly seasonal, and a trend modelled as a random walk.

This corresponds to a model for Poisson counts

$$y_k \sim Po(\mu_k),$$



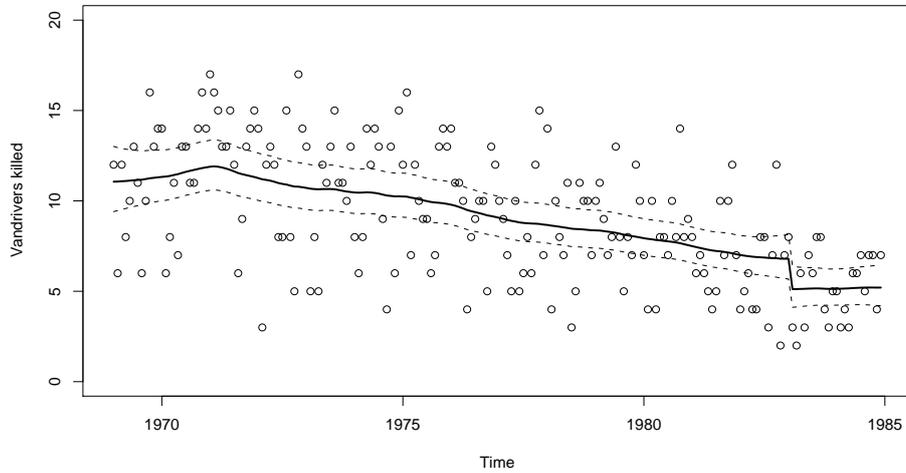


Figure 3: Estimated trend and intervention (solid line) for the vandriers data. The dashed lines are  $\pm 2$  standard deviations.

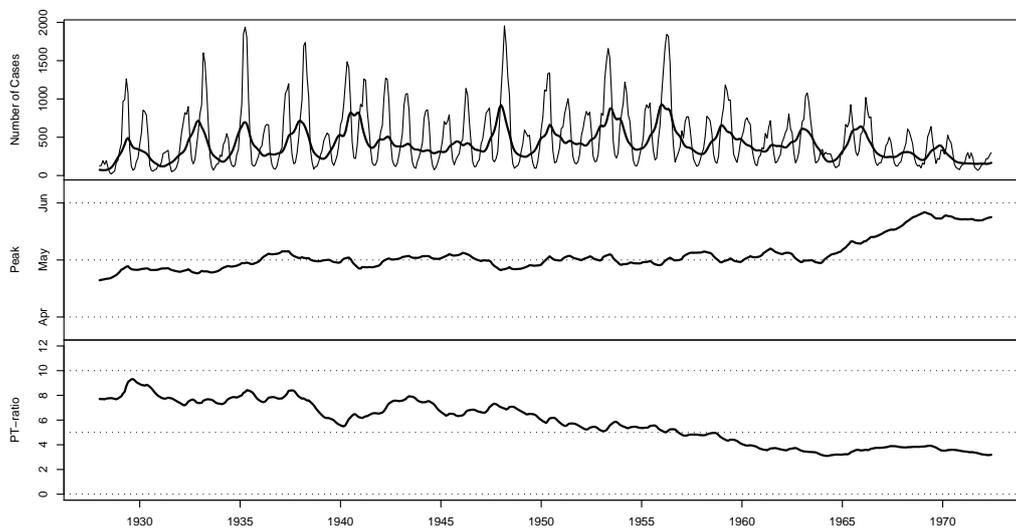


Figure 4: The variation in the incidence in mumps, NYC, 1927 – 1972. The upper frame shows the observed number of cases with the de-seasonalized trend superimposed. The middle frame shows the location of the peak of the seasonal pattern. The lower frame depicts the variation in the peak-to-trough ratio over the period.

Fitting a Poisson generalized linear model with a quadratic trend and an monthly seasonal pattern, yields an overdispersion of 89.7, a significant trend and a significant seasonal variation. Changing the seasonal pattern to a harmonic pattern is in accordance with the data but does not substantially change the overdispersion.

We model the mumps incidence with a first order polynomial trend with time-varying coefficients and a time-varying harmonic seasonal component.

$$y_k \sim \text{Po}(\mu_k),$$

with a linear trend and a harmonic seasonal, yielding the linear predictor

$$\lambda_k = \log \mu_k = T_k + H_k,$$

where the linear trend component is defined as

$$T_k = T_{k-1} + \beta_{k-1} + \omega_k^{(1)}, \quad \beta_k = \beta_{k-1} + \omega_k^{(2)},$$

and the harmonic seasonal pattern as

$$H_k = \theta_{ck} \cos\left(\frac{2\pi}{12}k\right) + \theta_{sk} \sin\left(\frac{2\pi}{12}k\right),$$

$$\theta_{ck} = \theta_{c,k-1} + \omega_k^{(c)}, \quad \theta_{sk} = \theta_{s,k-1} + \omega_k^{(s)},$$

$\omega_k^{(1)}$ ,  $\omega_k^{(2)}$ ,  $\omega_k^{(c)}$  and  $\omega_k^{(s)}$  being independent noise terms. Letting  $\boldsymbol{\theta}_k = [T_k \quad \beta_k \quad \theta_k^{(c)} \quad \theta_k^{(s)}]^\top$  we get the following matrix notation

$$\lambda_k = \begin{bmatrix} 1 & 0 & \cos\left(\frac{2\pi}{12}k\right) & \sin\left(\frac{2\pi}{12}k\right) \end{bmatrix} \boldsymbol{\theta}_{k-1},$$

and

$$\boldsymbol{\theta}_k = \begin{bmatrix} 1 & 1 & & \\ 0 & 1 & & \\ & & 1 & 0 \\ & & 0 & 1 \end{bmatrix} \boldsymbol{\theta}_{k-1} + \begin{bmatrix} \omega_k^{(1)} \\ \omega_k^{(2)} \\ \omega_k^{(c)} \\ \omega_k^{(s)} \end{bmatrix}.$$

This is formulated in **sspir** by the call

```
data("mumps")
index <- 1:length(mumps)
mumps.m <- ssm( mumps ~ -1 + tvar(polytime(index,1)) +
               tvar(polytrig(index,12,1)), family=poisson(link=log),
               phi=c(0,0,0.0005,0.0001),
               CO = diag(4))
mumps.fit <- getFit(mumps.m)

plot(mumps)
lines( exp(mumps.fit$m[,1]), lwd=2)
```

The choice of a first order sinusoid gives the possibility to express the seasonal variation via the peak-to-trough ratio (yearly max/min) and the location of the peak (code not shown, but available in `demo("mumps")`). The output in Figure 4 shows a gradually changing seasonal

pattern with a decreasing peak-to-trough ratio and a peak location slowly changing. The location of the peak is changing from late April in 1928–1936, where after the location of the peak stabilizes around May 1st until 1964, when the peak wanders off to late May, see Figure 4. It is also seen that the peak-to-trough ratio is varying between 6 to 9 until around 1948, when the ratio gradually decreases to about 4 in 1971. Epidemic episodes are seen irregularly each 3 to 5 years. The analysis can be reproduced in **sspir** by `demo("mumps")`.

## 5. Discussion

The main contribution of the **sspir** package is to give a formula language for specifying dynamic generalized linear models. That is, an extension of **glm** formulae by marking terms with `tvar` to specify that the corresponding coefficients are time-varying. The package also provides (extended) Kalman filter and Kalman smoother for models within the Gaussian, Poisson and binomial families. The output from the Kalman smoother leaves many possibilities for designing a suitable presentation of features of the latent process.

The Kalman filter is initialized by the values of  $\mathbf{m}_0$  and  $\mathbf{C}_0$ , see (3). The modeller can set entries in  $\mathbf{C}_0$  to accommodate prior knowledge. In cases where the prior information about  $\boldsymbol{\theta}_0$  is sparse, a diffuse initialization may be adequate, see Durbin and Koopman (2001). This feature has not yet been implemented.

The present framework does not allow the modeller to estimate the unknown variance parameters automatically. The modeller can, though, combine numerical maximization algorithms with the output of the iterated extended Kalman smoother. Hence, the formulation in **sspir** does not rely on any specific implementation of an inferential procedure.

## Acknowledgements

The authors are grateful for the helpful comments from Stefan Christensen, Spencer Graves, and two anonymous referees.

## References

- Diggle PJ, Heagerty PJ, Liang KY, Zeger SL (2002). *Analysis of Longitudinal Data*. Oxford University Press, 2nd edition.
- Durbin J, Koopman SJ (1997). “Monte Carlo Maximum Likelihood Estimation for Non-Gaussian State Space Models.” *Biometrika*, **84**(3), 669–684.
- Durbin J, Koopman SJ (2000). “Time Series Analysis of Non-Gaussian Observations Based on State Space Models from both Classical and Bayesian Perspectives.” *Journal of the Royal Statistical Society B*, **62**(1), 3–56. With discussion.
- Durbin J, Koopman SJ (2001). *Time Series Analysis by State Space Methods*. Oxford University Press.

- Gamerman D (1998). “Markov Chain Monte Carlo for Dynamic Generalised Linear Models.” *Biometrika*, **85**, 215–227.
- Harvey AC, Durbin J (1986). “The Effects of Seat Belt Legislation on British Road Casualties: A Case Study in Structural Time Series Modelling.” *Journal of the Royal Statistical Society A*, **149**(3), 187–227. With discussion.
- Hipel KW, McLeod IA (1994). *Time Series Modeling of Water Resources and Environmental Systems*. Elsevier Science Publishers B.V. (North-Holland).
- Kitagawa G, Gersch W (1984). “A Smoothness Priors-State Space Modeling of Time Series with Trend and Seasonality.” *Journal of the American Statistical Association*, **79**(386), 378–389.
- Koopman SJ, Shephard N, Doornik JA (1999). “Statistical Algorithms for Models in State Space Using **SsfPack** 2.2.” *Econometrics Journal*, **2**, 113–166.
- McCullagh P, Nelder JA (1989). *Generalized Linear Models*. Chapman and Hall, London, 2nd edition.
- R Development Core Team (2006). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org/>.
- Ripley BD (2002). “Time Series in R 1.5.0.” *R News*, **2**(2), 2–7.
- West M, Harrison PJ, Migon HS (1985). “Dynamic Generalized Linear Models and Bayesian Forecasting.” *Journal of the American Statistical Association*, **80**, 73–97. With discussion.
- Zeger SL (1988). “A Regression Model for Time Series of Counts.” *Biometrika*, **75**(4), 621–629.

**Affiliation:**

Claus Dethlefsen  
Center for Cardiovascular Research  
Aalborg Hospital, Aarhus University Hospital  
9000 Aalborg, Denmark  
E-mail: [aas.claus.dethlefsen@nja.dk](mailto:aas.claus.dethlefsen@nja.dk)  
URL: <http://www.math.aau.dk/~dethlef/>