

In this document we give a summary of the DLM time-series models implemented in this code package. We will describe a DLM model for atmospheric time-series with four key components: a smoothly varying non-linear background trend, a seasonal cycle comprised of two components with 12- and 6-month periods respectively, forcings described by a number regressor variables, and an auto-regressive process. We follow closely the notation used in [Laine, Latva-Pukkila & Kyrölä \(2014\)](#). Whilst we aim to give a clear and concise description of the model(s), this document is not meant to serve as a pedagogical introduction to DLM regression. For a more comprehensive review of DLM time-series analysis, see [Durbin & Koopman \(2012\)](#).

## DLM MODEL

In this section we describe the general DLM model, of which the various models implemented in this code package are special cases (described later on).

We model time-series observations  $y_t$  at times  $t$  as follows:

$$\begin{aligned} y_t &= \mathbf{F}_t^T \mathbf{x}_t + v_t, \quad v_t \sim \mathcal{N}(\mathbf{0}, V_t), \\ \mathbf{x}_t &= \mathbf{G}_t \mathbf{x}_{t-1} + \mathbf{w}_t, \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{W}_t), \end{aligned} \quad (1)$$

where the  $v_t$  is assumed zero-mean Gaussian measurement error with (time-dependent) variance  $V_t$ ,  $\mathbf{x}_t$  is the *state* containing the (non-linear) trend, the seasonal cycle, dynamical coefficients of the regressor variables, and the auto-regressive (AR) process, and the unobserved state  $\mathbf{x}_t$  is projected onto the observations through the *observation matrix*  $\mathbf{F}_t$ . The stochastic evolution of the state is governed by the *state evolution operator*  $\mathbf{G}_t$  and the *model error covariance*  $\mathbf{W}_t$ .

We will assume the state is constituted by the (non-linear) trend, seasonal cycle with two components (with 6- and 12-month periods respectively), a number of regressor variables and an ARn process. The state, observation, evolution and model covariance matrices hence decompose into these four components:

$$\begin{aligned} \mathbf{x}_t &= (\mathbf{x}_t^{\text{trend}}, \mathbf{x}_t^{\text{seas}}, \mathbf{x}_t^{\text{regressors}}, \mathbf{x}_t^{\text{AR}}), \quad \mathbf{F}_t = (\mathbf{F}_t^{\text{trend}}, \mathbf{F}_t^{\text{seas}}, \mathbf{F}_t^{\text{regressors}}, \mathbf{F}_t^{\text{AR}}), \\ \mathbf{G}_t &= \text{diag}(\mathbf{G}_t^{\text{trend}}, \mathbf{G}_t^{\text{seas}}, \mathbf{G}_t^{\text{regressors}}, \mathbf{G}_t^{\text{AR}}), \quad \mathbf{W}_t = \text{diag}(\mathbf{W}_t^{\text{trend}}, \mathbf{W}_t^{\text{seas}}, \mathbf{W}_t^{\text{regressors}}, \mathbf{W}_t^{\text{AR}}) \end{aligned} \quad (2)$$

Let us describe each of these components in turn.

### Non-linear trend

We model the non-linear background trend with two hidden states,  $\mathbf{x}_t^{\text{trend}} = (\mu_t, \alpha_t)$ , where  $\mu_t$  is the trend value at time  $t$  and  $\alpha_t$  is the change from time  $t$  to  $t + 1$ , governed by

$$\begin{aligned} \mu_t &= \mu_{t-1} + \alpha_{t-1} \\ \alpha_t &= \alpha_{t-1} + \epsilon_{\text{trend}}, \end{aligned} \quad (3)$$

where  $\epsilon_{\text{trend}} \sim \mathcal{N}(0, \sigma_{\text{trend}})$  is a zero-mean Gaussian random variate with standard deviation  $\sigma_{\text{trend}}$ . Hence, the background trend is modeled by a locally linear model with slope  $\alpha_t$ , where the slope is allowed to vary stochastically in time.

The rate at which the slope can evolve, and hence the non-linearity of the background trend, is governed entirely by the parameter  $\sigma_{\text{trend}}$ ; larger values of  $\sigma_{\text{trend}}$  will lead to more non-linear background evolution (on shorter timescales), while the limit  $\sigma_{\text{trend}} \rightarrow 0$  recovers a linear trend model (with constant slope). In the fits we keep  $\sigma_{\text{trend}}$  as a free parameter to be fit simultaneously with the rest of the regression model; this way, the data can decide (within reasonable prior constraints) how non-linear the background trend should be.

The state, observation, evolution and model covariance matrices governing the non-linear trend part of the model are given by:

$$\mathbf{x}_t^{\text{trend}} = (\mu_t, \alpha_t), \quad \mathbf{F}_t^{\text{trend}} = (1, 0), \quad \mathbf{G}_t^{\text{trend}} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{W}_t^{\text{trend}} = \begin{pmatrix} 0 & 0 \\ 0 & \sigma_{\text{trend}}^2 \end{pmatrix} \quad (4)$$

The observation matrix ensures that only the trend  $\mu_t$  enters into the observations  $y_t$  (Eq. 1), and the evolution and model error covariance matrices ensure the trend evolves according to Eq. 3.

## Seasonal cycle

The two-component seasonal cycle is modeled as follows:

$$\mathbf{x}_t^{\text{seas}} = (u_{12,t}, u_{12,t}^*, u_{6,t}, u_{6,t}^*), \quad \mathbf{F}_t^{\text{seas}} = (1, 0, 1, 0), \quad \mathbf{G}_t^{\text{seas}} = \text{diag}(\mathbf{G}_{12,t}^{\text{seas}}, \mathbf{G}_{6,t}^{\text{seas}}), \quad \mathbf{W}_t^{\text{seas}} = \text{diag}(\mathbf{W}_{12,t}^{\text{seas}}, \mathbf{W}_{6,t}^{\text{seas}})$$

$$\mathbf{G}_{k,t}^{\text{seas}} = \begin{pmatrix} \cos(k 2\pi/12) & \sin(k 2\pi/12) \\ -\sin(k 2\pi/12) & \cos(k 2\pi/12) \end{pmatrix}, \quad \mathbf{W}_{k,t}^{\text{seas}} = \begin{pmatrix} \sigma_{\text{seas}}^2 & 0 \\ 0 & \sigma_{\text{seas}}^2 \end{pmatrix}, \quad (5)$$

where  $u_{k,t}$  denote the amplitudes of the orthogonal sinusoidal seasonal components, which can evolve stochastically in time. The degree to which the seasonal cycle amplitude and phase can evolve in time is governed entirely by the parameter  $\sigma_{\text{seas}}$ : setting  $\sigma_{\text{seas}} = 0$  corresponds to a constant amplitude and phase seasonal cycle, whilst  $\sigma_{\text{seas}} > 0$  allows for evolution in the amplitude and phase of the seasonal cycle over time. We keep  $\sigma_{\text{seas}}$  as an additional free parameter to be fit alongside the rest of the regression model, allowing for modulation of the seasonal variability if this is preferred by the data.

For monthly data, the 6- and 12-month seasonal components are described by setting  $k = 2$  and  $k = 1$  respectively. For daily or annual data these are scaled accordingly.

## Forcing via regressor variables

The forcing described by  $n$  regressor variables, where we have a corresponding time-series  $z_t$  for each regressor, enters into the DLM model as follows:

$$\mathbf{x}_t^{\text{regressors}} = (\beta_{1,t}, \beta_{2,t}, \dots, \beta_{n,t}), \quad \mathbf{F}_t^{\text{regressors}} = (z_{1,t}, z_{2,t}, \dots, z_{n,t}), \quad \mathbf{G}_t^{\text{regressors}} = \mathbf{I}_{n \times n},$$

$$\mathbf{W}_t^{\text{regressors}} = \text{diag}(\sigma_{z_1}^2, \sigma_{z_2}^2, \dots, \sigma_{z_n}^2), \quad (6)$$

where  $\mathbf{I}_{n \times n}$  is the identity matrix, and  $\sigma_{z_i}$  allows for time-dependent evolution of the amplitudes of the regressors  $\beta_{i,t}$ . The parameters  $\sigma_{z_i}$  controls the allowed degree of time-variability in the amplitude of the regressors, where  $\sigma_{z_i} = 0$  denotes constant (in time) amplitudes. In this case, the amplitudes  $\beta_i$  can be interpreted in much the same manner as for multiple linear regression analyses. However, in general  $\sigma_{z_i}$  can be kept as free parameters and fit alongside the rest of the regression model, so the data can determine whether time-varying regressor amplitudes are preferred.

## ARn process

The ARn process  $q_t \sim \text{AR}n$  enters into the DLM model as follows:

$$\mathbf{x}_t^{\text{AR}} = (q_t, q_{t-1}, \dots, q_{t-(n-1)}), \quad \mathbf{F}_t^{\text{AR}} = (1, 0, \dots, 0), \quad \mathbf{G}_t^{\text{AR}} = \begin{pmatrix} \rho_1 & \rho_2 & \dots & \rho_n \\ 1 & 0 & \dots & 0 \\ \vdots & & \ddots & 0 \\ 1 & 0 & \dots & 0 \end{pmatrix}, \quad \mathbf{W}_t^{\text{AR}} = \begin{pmatrix} \sigma_{\text{AR}}^2 & 0 & & \\ 0 & 0 & \dots & 0 \\ \vdots & & \ddots & \\ 0 & & & 0 \end{pmatrix}, \quad (7)$$

where the ARn parameters  $\rho_1, \rho_2, \dots, \rho_n$  and  $\sigma_{\text{AR}}$  are kept as free parameters that are fit with the regression model. In this code package we include models with AR1 and AR2 processes, but the code is easily extendible to higher order AR processes.

## MODEL PARAMETERS AND PRIORS

In order to perform Bayesian inference under the DLM model described above, we must define priors over the hyper-parameters and the initial state  $\mathbf{x}_0$ . We choose the following priors as standard in the code:

$$\begin{aligned} \sigma_{\text{trend}} &\sim \text{HalfNormal}(0, \sigma_{\text{trend}}^{\text{prior}}) \\ \sigma_{\text{seas}} &\sim \text{HalfNormal}(0, \sigma_{\text{seas}}^{\text{prior}}) \\ \sigma_z &\sim \text{HalfNormal}(\mathbf{0}, \sigma_z^{\text{prior}}) \\ \sigma_{\text{AR}} &\sim \text{HalfNormal}(0, \sigma_{\text{AR}}^{\text{prior}}) \\ \rho_i &\sim \text{Uniform}(-1, 1) \\ \mathbf{x}_0 &\sim \text{Normal}(\mathbf{0}, \mathbf{I} \times S). \end{aligned} \quad (8)$$

The user can control the widths of all priors (with exception of the AR correlations  $\rho_i$ ), but default values are taken to be:  $\sigma_{\text{trend}}^{\text{prior}} = 10^{-4}$ ,  $\sigma_{\text{seas}}^{\text{prior}} = 10^{-2}$ ,  $\sigma_z^{\text{prior}} = 10^{-4}\mathbf{I}$ ,  $S = 10$ .

Table 1: Description of all of the inferred parameters available to each specific model implementation and how to refer to those parameters in the code.  $N$  refers to the number of data points in the analyzed time-series,  $n_{\text{reg}}$  refers to the number of regressors that are used in the analysis.

| variable                | name in code | type                           | d1m_vanilla_AR1 | d1m_vanilla_AR2 | d1m_noregs_AR1 | d1m_dynregs_AR1 |
|-------------------------|--------------|--------------------------------|-----------------|-----------------|----------------|-----------------|
| $\mu_t$                 | trend        | np.array(N)                    | X               | X               | X              | X               |
| $\alpha_t$              | slope        | np.array(N)                    | X               | X               | X              | X               |
| $\beta_t$               | beta         | np.array(N, $n_{\text{reg}}$ ) | X               | X               |                | X               |
| $u_{12,t} + u_{6,t}$    | seasonal     | np.array(N)                    | X               | X               | X              | X               |
| $q_t$                   | ar           | np.array(N)                    | X               | X               | X              | X               |
| $\sigma_{\text{trend}}$ | sigma_trend  | float                          | X               | X               | X              | X               |
| $\sigma_{\text{seas}}$  | sigma_seas   | float                          | X               | X               | X              | X               |
| $\sigma_{\text{AR}}$    | sigma_AR     | float                          | X               | X               | X              | X               |
| $\sigma_z$              | sigma_reg    | np.array( $n_{\text{reg}}$ )   |                 |                 |                | X               |
| $\rho_1$                | rhoAR1       | float                          | X               | X               | X              | X               |
| $\rho_2$                | rhoAR2       | float                          |                 | X               |                |                 |

## SPECIFIC MODEL DESCRIPTIONS

In this package we include implementations of a number of DLM models that all inherit from the general DLM model described above, but with certain specific specifications (eg., with certain features turned on or off). We describe the specs of each implemented model below:

### d1m\_vanilla\_AR1

This is the go to model for most applications and has the following simplifications:  $\sigma_{z,i} = 0$  for all regressors (so the regressor amplitudes are assumed to be constant in time), and as the name suggest an AR1 auto-regressive process. This model was used for analysis of stratospheric ozone data in [Ball et al. \(2017\)](#) and [Ball et al. \(2018\)](#).

### d1m\_vanilla\_AR2

As above (constant in time regressor amplitudes) but with an AR2 process.

### d1m\_noregs\_AR1

This model has no regressors and is comprised only of a trend, seasonal cycle and auto-regressive (AR1) process.

### d1m\_dynregs\_AR1

This is the full general model described above, with an AR1 auto-regressive process.

## REFERENCES

- Ball W. T., Alsing J., Mortlock D. J., Rozanov E. V., Tummon F., Haigh J. D., 2017, Atmospheric Chemistry and Physics, 17, 12269
- Ball W. T. et al., 2018, Atmospheric Chemistry and Physics, 18, 1379
- Durbin J., Koopman S. J., 2012, Time series analysis by state space methods, Vol. 38. Oxford University Press
- Laine M., Latva-Pukkila N., Kyrölä E., 2014, Atmospheric Chemistry and Physics, 14, 9707