

# Linear State-Space Models

September 3, 2024

The material in these notes draws from several sources, including Chapter 6 of S&S, Chapter 11 of Chan (2010), and the online textbook [Applied Time Series Analysis](#) which was written by the authors of the MARSS software for R.

## Introduction to a Simple State-Space Model

State space models give us yet another way of writing out a model for our time series. In this section, the notation we use will differ slightly. The observed time series that we intend to analyze has values  $y_1, \dots, y_n$ . Let's work through a small example! For a univariate time series of length  $n$  we assume

$$y_t = ax_t + v_t \quad \text{Observation Equation}$$

$$x_t = \phi x_{t-1} + w_t \quad \text{State Equation,}$$

where  $v_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_v^2)$ ,  $w_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_w^2)$  and  $x_1 = \mu$ . The state equation looks like an **AR**(1) model, but values of the **AR**(1) process  $x_t$  are not observed, and it differs from an **AR**(1) model insofar as the initial state is assumed to be fixed. Rather, a linear transformation  $y_t$  of the states  $x_t$  are observed with some additional normal noise. To emphasize, we observe values of  $y_t$ , but we do *not* observe values of  $x_t$ .

## Predicting, Filtering, and Smoothing the States and Forecasting the Time Series

For now, let's assume that we know  $\sigma_v^2$ ,  $\sigma_w^2$ ,  $a$ ,  $\phi$ , and  $\mu$ . Having assumed a state-space model, we might be interested in:

- **Predicting**, i.e. estimating future values of  $x_t$  given past values of  $y_{t-1}, \dots, y_1$ 
  - For example, computing  $\mathbb{E}[x_2|y_1]$  and  $\mathbb{V}[x_2|y_1]$
- **Filtering**, i.e. estimating future values of  $x_t$  given current and past values  $y_t, \dots, y_1$ 
  - For example, computing  $\mathbb{E}[x_2|y_1, y_2]$  and  $\mathbb{V}[x_2|y_1, y_2]$
- **Smoothing**, i.e. estimating a past value of  $x_{t-1}$  given all observed past, current, and future values  $y_n, \dots, y_1$ 
  - For example, computing  $\mathbb{E}[x_1|y_1, y_2]$  and  $\mathbb{V}[x_1|y_1, y_2]$

As usual, we might also be interested in forecasting future values  $y_{n+k}$  given  $y_1, \dots, y_n$ , which will require **prediction** of the states because we will need to compute  $\mathbb{E}[x_{n+k}|y_1, \dots, y_n]$ . Even if we do know the values of these parameters, how do we go about predicting, filtering, and smoothing?

The first step is to realize that we can write out the joint probability distribution of  $\mathbf{y}$  and  $\mathbf{x}$ . Because we have assumed that the errors  $\mathbf{v}$  and  $\mathbf{w}$  and the initial value  $x_1$  is fixed, we know that the joint distribution of  $\mathbf{y}$  and  $\mathbf{x}$  is given by:

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{x} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mathbb{E}[\mathbf{y}] \\ \mathbb{E}[\mathbf{x}] \end{pmatrix}, \begin{pmatrix} \mathbb{V}[\mathbf{y}] & \text{Cov}[\mathbf{y}, \mathbf{x}] \\ \text{Cov}[\mathbf{y}, \mathbf{x}] & \mathbb{V}[\mathbf{x}] \end{pmatrix} \right). \quad (1)$$

Why is this useful? Well, predicting, filtering, smoothing, and forecasting all correspond to different conditional mean and variance estimation problems, and when we have a normal

joint distribution and know its parameters, it's easy to figure out any conditional distribution we want! We can see this most easily in the case of smoothing, which corresponds to the problem of estimating the conditional means  $\mathbb{E}[\mathbf{x}|\mathbf{y}]$  and variances  $\mathbb{V}[\mathbf{x}|\mathbf{y}]$ . Before we go further, a very useful property of the multivariate normal distribution is that:

$$\text{If } \begin{pmatrix} \mathbf{u} \\ \mathbf{v} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}, \begin{pmatrix} \mathbf{C} & \mathbf{D} \\ \mathbf{D}' & \mathbf{E} \end{pmatrix} \right), \text{ then}$$

– The marginal distributions of  $\mathbf{u}$  and  $\mathbf{v}$  are given by

$$\mathbf{u} \sim \mathcal{N}(\mathbf{a}, \mathbf{C}) \text{ and } \mathbf{v} \sim \mathcal{N}(\mathbf{b}, \mathbf{E})$$

– The conditional distributions of  $\mathbf{u}|\mathbf{v}$  and  $\mathbf{v}|\mathbf{u}$

$$\mathbf{u}|\mathbf{v} \sim \mathcal{N}(\mathbf{a} + \mathbf{D}'\mathbf{E}^{-1}(\mathbf{v} - \mathbf{b}), \mathbf{C} - \mathbf{D}'\mathbf{E}^{-1}\mathbf{D})$$

$$\mathbf{v}|\mathbf{u} \sim \mathcal{N}(\mathbf{b} + \mathbf{D}\mathbf{C}^{-1}(\mathbf{u} - \mathbf{a}), \mathbf{E} - \mathbf{D}\mathbf{C}^{-1}\mathbf{D}')$$

We can use these facts to find the conditional distributions of the states  $\mathbf{x}$  given the data  $\mathbf{y}$  and accordingly, the smoothed estimates of the states, but first we'll simplify the joint distribution a bit. We can simplify  $\mathbb{E}[\mathbf{y}]$ ,  $\mathbb{V}[\mathbf{y}]$ , and  $\text{Cov}[\mathbf{x}, \mathbf{y}]$ .

- It's pretty straightforward to see that  $\mathbb{E}[\mathbf{y}] = a\mathbb{E}[\mathbf{x}]$
- A bit of algebra lets us rewrite  $\mathbb{V}[\mathbf{y}]$

$$\begin{aligned} \mathbb{V}[\mathbf{y}] &= \mathbb{E}[(\mathbf{y} - \mathbb{E}[\mathbf{y}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])'] \\ &= \mathbb{E}[(a(\mathbf{x} - \mathbb{E}[\mathbf{x}]) + \mathbf{v})(a(\mathbf{x} - \mathbb{E}[\mathbf{x}]) + \mathbf{v})'] \\ &= a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n. \end{aligned}$$

- A bit more algebra lets us rewrite  $\text{Cov}[\mathbf{x}, \mathbf{y}]$

$$\begin{aligned}\text{Cov}[\mathbf{x}, \mathbf{y}] &= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])'] \\ &= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(a\mathbf{x} - a\mathbb{E}[\mathbf{x}] + \mathbf{w})'] \\ &= a\mathbb{V}[\mathbf{x}]\end{aligned}$$

Plugging these expressions into 2 yields a nicely structured normal distribution,

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{x} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} a\mathbb{E}[\mathbf{x}] \\ \mathbb{E}[\mathbf{x}] \end{pmatrix}, \begin{pmatrix} a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n & a\mathbb{V}[\mathbf{x}] \\ a\mathbb{V}[\mathbf{x}] & \mathbb{V}[\mathbf{x}] \end{pmatrix} \right). \quad (2)$$

Using our normal distribution facts and (2), we can return to the smoothing problem

$$\begin{aligned}\mathbf{x}|\mathbf{y} &\sim \mathcal{N}(\mathbb{E}[\mathbf{x}] + a\mathbb{V}[\mathbf{x}](a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n)^{-1}(\mathbf{y} - a\mathbb{E}[\mathbf{x}]), \\ &\quad \mathbb{V}[\mathbf{x}] - a^2\mathbb{V}[\mathbf{x}](a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n)^{-1}\mathbb{V}[\mathbf{x}]).\end{aligned}$$

The smoothed values of the states are given by the conditional mean

$$\mathbb{E}[\mathbf{x}] + a\mathbb{V}[\mathbf{x}](a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n)^{-1}(\mathbf{y} - a\mathbb{E}[\mathbf{x}]),$$

and our uncertainty about them is quantified by the conditional variance,

$$\mathbb{V}[\mathbf{x}] - a^2\mathbb{V}[\mathbf{x}](a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n)^{-1}\mathbb{V}[\mathbf{x}].$$

Predicting and filtering can be done similarly by identifying the marginal joint distribution of a state and the corresponding observed values, and computing the corresponding conditional distribution.

- We can obtain predicted means and variances of the states  $x_t$  given the past observed values of the time series  $y_1, \dots, y_{t-1}$  by letting  $u = x_t$  and  $\mathbf{v} = (y_1, \dots, y_{t-1})$ , using our marginal and conditional normal distribution facts;
- We can obtain filtered means and variances of the states  $x_t$  given the past and current

observed values of the time series  $y_1, \dots, y_t$  using our marginal and conditional normal distribution facts, letting  $u = x_t$  and  $\mathbf{v} = (y_1, \dots, y_t)$ .

Forecasts follow straightforwardly. We have assumed  $y_t = ax_t + v_t$ , so it is natural to define a forecasted value of  $y_{n+k}$  as

$$\mathbb{E}[y_{n+k}|y_1, \dots, y_n] = a\mathbb{E}[x_{n+k}|y_1, \dots, y_n],$$

with variance

$$\mathbb{V}[y_{n+k}|y_1, \dots, y_n] = a^2\mathbb{V}[x_{n+k}|y_1, \dots, y_n] + \sigma_w^2.$$

In general, even though we *can* figure out any conditional distribution we want, it would be very computationally burdensome to figure out each conditional expectation one by one. Fortunately, the **Kalman filter** allows for fast, recursive estimation of:

- Predicted means and variances
- Filtered means and variances
- Smoothed means and variances
- Forecast means and variances.

The existence of a fast algorithm to compute all of these quantities is essential in practice.

- **Kalman Prediction & Filtering:** The Kalman prediction and filtering algorithms work together, iteratively computing the predictions, prediction variances, filtered values of  $x_t$ , and variances of the filtered values, starting at time  $t = 1$ . To aid interpretation of the algorithm, predictions and their variances are printed in **blue**, while filtered values and their variances are printed in **red**.

$$\begin{aligned} - \mathbb{E}[x_t|y_1, \dots, y_{t-1}] &= \phi \mathbb{E}[x_{t-1}|y_1, \dots, y_{t-2}] + \left( \frac{a\phi\mathbb{V}[x_{t-1}|y_1, \dots, y_{t-2}]}{a^2\mathbb{V}[x_{t-1}|y_1, \dots, y_{t-2}] + \sigma_v^2} \right) (y_{t-1} - a\mathbb{E}[x_{t-1}|y_1, \dots, y_{t-2}]) \\ - \mathbb{V}[x_t|y_1, \dots, y_{t-1}] &= \sigma_w^2 + \phi^2\mathbb{V}[x_{t-1}|y_1, \dots, y_{t-2}] \left( 1 - \frac{a^2\mathbb{V}[x_{t-1}|y_1, \dots, y_{t-2}]}{a^2\mathbb{V}[x_{t-1}|y_1, \dots, y_{t-2}] + \sigma_v^2} \right) \\ - \mathbb{E}[x_t|y_1, \dots, y_t] &= \mathbb{E}[x_t|y_1, \dots, y_{t-1}] + \left( \frac{a\mathbb{V}[x_t|y_1, \dots, y_{t-1}]}{a^2\mathbb{V}[x_t|y_1, \dots, y_{t-1}] + \sigma_v^2} \right) (y_t - a\mathbb{E}[x_t|y_1, \dots, y_{t-1}]) \end{aligned}$$

$$- \mathbb{V}[x_t|y_1, \dots, y_t] = \mathbb{V}[x_t|y_1, \dots, y_{t-1}] \left( 1 - \frac{a^2 \mathbb{V}[x_t|y_1, \dots, y_{t-1}]}{a^2 \mathbb{V}[x_t|y_1, \dots, y_{t-1}] + \sigma_v^2} \right)$$

- **Kalman Smoothing:** Kalman smoothing works in reverse, starting at time  $t = n - 1$ . It uses the Kalman predictions, filtered values, and their variances. For each value of  $t$ , the smoother starts at  $s = 1$  and proceeds until  $s = n - t$ .

- The smoother first computes an expectation:

$$\mathbb{E}[x_t|y_1, \dots, y_{t+s}] = \mathbb{E}[x_t|y_1, \dots, y_{t+s-1}] + \frac{a \mathbb{E}[(x_t - \mathbb{E}[x_t|y_1, \dots, y_{t-1}]) (x_{t+s} - \mathbb{E}[x_{t+s}|y_1, \dots, y_{t+s-1}])]}{a^2 \mathbb{V}[x_{t+s}|y_1, \dots, y_{t+s-1}] + \sigma_v^2} (y_{t+s} - a \mathbb{E}[x_{t+s}|y_1, \dots, y_{t+s-1}])$$

- Then some covariances as well:

$$\mathbb{E}[(x_t - \mathbb{E}[x_t|y_1, \dots, y_{t-1}]) (x_{t+s} - \mathbb{E}[x_{t+s}|y_1, \dots, y_{t+s-1}])] = \phi \mathbb{E}[(x_t - \mathbb{E}[x_t|y_1, \dots, y_{t-1}]) (x_{t+s-1} - \phi \mathbb{E}[x_{t+s-1}|y_1, \dots, y_{t+s-2}])] \left( 1 - \frac{a^2 \mathbb{V}[x_{t+s-1}|y_1, \dots, y_{t+s-2}]}{a^2 \mathbb{V}[x_{t+s-1}|y_1, \dots, y_{t+s-2}] + \sigma_v^2} \right)$$

- And finally a variance:

$$\mathbb{V}[x_t|y_1, \dots, y_{t+s}] = \mathbb{V}[x_t|y_1, \dots, y_{t+s-1}] - a^2 \left( \frac{\mathbb{E}[(x_t - \mathbb{E}[x_t|y_1, \dots, y_{t-1}]) (x_{t+s} - \mathbb{E}[x_{t+s}|y_1, \dots, y_{t+s-1}])]^2}{a^2 \mathbb{V}[x_{t+s}|y_1, \dots, y_{t+s-1}] + \sigma_v^2} \right)$$

The intuition motivating these algorithms is closely related to the intuition behind the Durbin-Levinson algorithm, which we saw when we studied forecasting. The basic idea is that the algorithms let us avoid repeated computationally expensive matrix inversions by exploiting relationships between the values we want to compute.

## Estimating the State-Space Model Parameters

So far, we've assumed that the parameters  $\sigma_v^2$ ,  $\sigma_w^2$ ,  $a$ ,  $\phi$ , and  $\mu$  are known. In practice,  $a$  is often assumed to be fixed and known at  $a = 1$  but we may need to estimate the rest of the parameters from the data. This requires maximizing the **marginal likelihood** of the data  $\mathbf{y}$ , having integrated the latent time series  $\mathbf{x}$  out. This is given by:

$$p(\mathbf{y}|\sigma_v^2, \sigma_w^2, a, \phi, \mu) = \int p(\mathbf{y}|\mathbf{x}, \sigma_v^2, a) p(\mathbf{x}|\mu, \phi, \sigma_w^2) d\mathbf{x}. \quad (3)$$

Maximizing over an integral is tricky!

## Direct Maximum Marginal Likelihood

Fortunately, our normal distribution facts tell us that the marginal distribution of  $\mathbf{y}$  is

$$\mathbf{y} \sim \mathcal{N}(a\mathbb{E}[\mathbf{x}], a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n), \quad (4)$$

which means that the integral can be evaluated to a normal density for  $\mathbf{y}$  equal to:

$$\frac{1}{\sqrt{2\pi}^n \sqrt{|a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n|}} \exp \left\{ -(\mathbf{y} - a\mathbb{E}[\mathbf{x}])' (a^2\mathbb{V}[\mathbf{x}] + \sigma_v^2\mathbf{I}_n)^{-1} (\mathbf{y} - a\mathbb{E}[\mathbf{x}]) \right\},$$

where  $|\mathbf{A}|$  refers to the determinant of the matrix  $\mathbf{A}$ . If we could maximize this over  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$ , then we'd be all set! Unfortunately, evaluating the marginal likelihood requires inverting an  $n \times n$  matrix, which can be very computationally expensive. It would be better if we could write the marginal likelihood as a function of just the data and output from the Kalman predictor, filter, or smoother. Fortunately, this is possible!

Let  $r_t = y_t - a\mathbb{E}[x_t|y_1, \dots, y_{t-1}]$  be the one-step-ahead forecast errors with variances  $\mathbb{V}[r_t|y_1, \dots, y_{t-1}] = a^2\mathbb{V}[x_t|y_1, \dots, y_{t-1}] + \sigma_v^2$ . Conveniently, it's possible to use our normal distribution facts to show that we can write the likelihood in terms of the residuals as

$$p(r_1) \prod_{t=2}^n p(r_t|y_1, \dots, y_{t-1}) = \frac{1}{\sqrt{2\pi}^n} \frac{1}{\sqrt{\mathbb{V}[r_1]}} \exp \left\{ -\frac{r_1^2}{2\mathbb{V}[r_1]} \right\} \prod_{t=2}^n \frac{1}{\sqrt{\mathbb{V}[r_t|y_1, \dots, y_{t-1}]}} \exp \left\{ -\frac{r_t^2}{2\mathbb{V}[r_t|y_1, \dots, y_{t-1}]} \right\}. \quad (5)$$

Defining the likelihood in this way doesn't quite get us out of the woods - although we can evaluate it quickly given values of  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$ , it is still difficult to maximize (5) over  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  jointly as written because they all affect the values of  $\mathbf{r}$  as well as the variances. Fortunately, we can take a computationally simpler approach which maximizes (5) by picking starting values of the parameters  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  and then iteratively:

- Fixing  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  and computing  $\mathbf{r}$  accordingly, using the Kalman prediction algorithm;
- Maximizing (5) over  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  for fixed  $\mathbf{r}$ .

This iterative process continues until the likelihood (6) or successive values of the parameters

$\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  converge. This is a popular way to compute maximum likelihood estimates of  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$ , but it can be difficult to work in practice because the second step of maximizing (5) over  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  for fixed  $\mathbf{r}$  is still a tricky nonlinear optimization problem.

## Expectation-Maximization (EM) Maximum Marginal Likelihood

One other way we can compute maximum likelihood estimates of  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  is to use the expectation-maximization (EM) algorithm, which is an algorithm for maximizing a function that corresponds to an integral over some latent variables, which in this case are the latent states  $x_t$ . The EM algorithm allows us to maximize the marginal likelihood in (3) by maximizing  $\mathbb{E}[\log(p(\mathbf{y}|\mathbf{x}, \sigma_v^2, a) p(\mathbf{x}|\mu, \phi, \sigma_w^2)) | y_1, \dots, y_n]$  using an iterative procedure. It isn't immediately obvious how this helps us, but we will get a sense by working simplifying the conditional expectation of the joint log-likelihood:

$$\begin{aligned}
& \mathbb{E}[\log(p(\mathbf{y}|\mathbf{x}, \sigma_v^2, a) p(\mathbf{x}|\mu, \phi, \sigma_w^2)) | y_1, \dots, y_n] = \\
& \mathbb{E}\left[\log\left(p(y_1|x_1, \sigma_v^2, a) p(x_1|\mu) \prod_{t=2}^n p(y_t|x_t, \sigma_v^2, a) p(x_t|x_{t-1}, \phi, \sigma_w^2)\right) | y_1, \dots, y_n\right] = \\
& K - \frac{n}{2}\log(\sigma_v^2) - \frac{1}{2\sigma_v^2} \sum_{t=1}^n \mathbb{E}\left[(y_t - ax_t)^2 | y_1, \dots, y_n\right] - \\
& \frac{(n-1)}{2}\log(\sigma_w^2) - \frac{1}{2\sigma_w^2} \sum_{t=2}^n \mathbb{E}\left[(x_t - \phi x_{t-1})^2 | y_1, \dots, y_n\right] = \\
& K - \frac{n}{2}\log(\sigma_v^2) - \frac{1}{2\sigma_v^2} \sum_{t=1}^n \mathbb{E}\left(y_t^2 - 2ay_t\mathbb{E}[x_t|y_1, \dots, y_n] + a^2\mathbb{E}[x_t^2|y_1, \dots, y_n]\right) - \tag{6} \\
& \frac{(n-1)}{2}\log(\sigma_w^2) - \frac{1}{2\sigma_w^2} \sum_{t=3}^n \left(\mathbb{E}[x_t^2|y_1, \dots, y_n] - 2\phi\mathbb{E}[x_t x_{t-1}|y_1, \dots, y_n] + \phi^2\mathbb{E}[x_{t-1}^2|y_1, \dots, y_n]\right) - \\
& \frac{1}{2\sigma_w^2} \left(\mathbb{E}[x_2^2|y_1, \dots, y_n] - 2\phi\mu\mathbb{E}[x_2|y_1, \dots, y_n] + \phi^2\mu^2\right),
\end{aligned}$$

where  $x_1 = \mu$ . where  $K$  is a constant that doesn't depend on the data  $\mathbf{y}$ , or the latent states  $\mathbf{x}$ , or the state-space parameters  $a, \phi, \sigma_v^2, \sigma_w^2$ , and  $\mu$ . We can maximize the marginal likelihood in (3) by picking starting values of the parameters  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  and then iteratively:

- Fixing the parameter values  $a, \phi, \sigma_v^2, \sigma_w^2$ , and  $\mu$  and computing the expectations  $\mathbb{E}[x_t|y_1, \dots, y_n]$  and  $\mathbb{E}[x_t x_{t-1}|y_1, \dots, y_n]$  using output from the Kalman smoother;



- Fixing the expectations  $\mathbb{E}[x_t|y_1, \dots, y_n]$  and  $\mathbb{E}[x_t x_{t-1}|y_1, \dots, y_n]$  and maximizing (6) over  $a, \phi, \sigma_v^2, \sigma_w^2$ , and  $\mu$ .

This iterative process continues until the likelihood (6) or successive values of the parameters  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  converge. It's common to refer to the first step as the *E*-step, because it involves computing expectations, and it is common to refer to the second step as the *M*-step, because it involves maximizing a function. This approach is computationally much simpler than the previous one, because the parameters  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  enter into the *M*-step in much more convenient ways. In fact, there are closed form solutions to the *M*-step, i.e. we can write down formulas for the values of  $\sigma_v^2, \sigma_w^2, a, \phi$ , and  $\mu$  that maximize (6) when the expectations  $\mathbb{E}[x_t|y_1, \dots, y_n]$  and  $\mathbb{E}[x_t x_{t-1}|y_1, \dots, y_n]$  are held constant! The main disadvantage to using this approach is that it can converge very slowly. In practice, we often use a hybrid of both approaches, by using the results of a (possibly unconverged) EM algorithm as starting values for a direct maximum marginal likelihood estimation procedure.

## More General Linear State-Space Models

Let's make things a little more general! We can incorporate covariates  $\mathbf{z}_t$ , and allow the linear transformation from the states  $x_t$  to the observed data  $y_t$  to depend on time by letting  $a$  depend on time, denoted by  $a_t$ :

$$y_t = a_t x_t + \mathbf{z}_t' \boldsymbol{\gamma} + v_t \quad \text{Observation Equation}$$

$$x_t = \phi x_{t-1} + \mathbf{z}_t' \mathbf{v} + w_t \quad \text{State Equation,}$$

where  $v_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_v^2)$ ,  $w_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_w^2)$  and  $x_1 \sim \mathcal{N}(\mu, \sigma_x^2)$  and  $a_t$  is either fixed and known or known up to a constant, i.e.  $a_t = a c_t$  where  $c_t$  is fixed and known but  $a$  is not. The Kalman prediction, filtering, and smoothing algorithms we discussed in the simpler case can be extended to quickly compute predictions, filtered values, and smoother values for the more general model. The maximum likelihood estimation procedures can likewise be generalized to estimate  $\sigma_v^2, \sigma_w^2, \phi, a, \boldsymbol{\gamma}, \mathbf{v}$ , and  $\mu$  from the marginal likelihood of the observed time series.

In state-space models, **seasonal effects** are often included via the inclusion of covariates in the observed data process,  $\mathbf{z}_t$ . For example, if the data is observed daily but weekly effects are expected to be present, we might include day of the week indicators in  $\mathbf{z}_t$ . Alternatively, we could include the covariates  $\mathbf{z}_t$  in the state process. Usually, we only include the covariates in one process or the other, i.e. we either set  $\boldsymbol{\gamma} = \mathbf{0}$  and estimate  $\mathbf{v}$  from the data or we set  $\mathbf{v} = \mathbf{0}$  and estimate  $\boldsymbol{\gamma}$  from the data. Additionally, state-space models can be used to perform **stochastic regression** with a single covariate, the values of  $a_t$  correspond to the values of a single covariate. If this is the case, we can think of  $x_t$  as a time-varying regression coefficient.

### Getting Into the Weeds

It can be valuable to derive closed form expressions for  $\mathbb{E}[\mathbf{x}]$  and  $\mathbb{V}[\mathbf{x}]$  to delve into the implied marginal distribution of the data under the state-space models discussed here. First, we derive  $\mathbb{E}[\mathbf{x}]$  and  $\mathbb{V}[\mathbf{x}]$ .

$$\begin{aligned}
x_1 &= \mu + \mathbf{z}'_1 \mathbf{v} \\
x_2 &= \phi(\mu + \mathbf{z}'_1 \mathbf{v}) + \mathbf{z}'_2 \mathbf{v} + w_2 \\
x_3 &= \phi x_2 + w_3 = \phi^2(\mu + \mathbf{z}'_1 \mathbf{v}) + \phi \mathbf{z}'_2 \mathbf{v} + \mathbf{z}'_3 \mathbf{v} + \phi w_2 + w_3 \\
x_4 &= \phi x_3 + w_4 = \phi^3(\mu + \mathbf{z}'_1 \mathbf{v}) + \phi^2 \mathbf{z}'_2 \mathbf{v} + \phi \mathbf{z}'_3 \mathbf{v} + \mathbf{z}'_4 \mathbf{v} + \phi^2 w_2 + \phi w_3 + w_4 \\
x_k &= \phi^{k-1} \mu + \left( \sum_{i=1}^k \phi^{k-i} \mathbf{z}'_i \right) \mathbf{v} + \left( \sum_{i=2}^k \phi^{k-i} w_i \right) \\
\mathbb{E}[x_t] &= \phi^{t-1} \mu + \left( \sum_{i=1}^k \phi^{k-i} \mathbf{z}'_i \right) \mathbf{v} \\
\mathbb{V}[x_t] &= \sigma_w^2 \sum_{i=2}^t \phi^{2(t-i)} \\
\text{Cov}[x_t, x_s] &= \sigma_w^2 \sum_{i=2}^{\min\{s,t\}} \phi^{t+s-2i}, \quad s, t > 1 \\
\text{Cov}[x_1, x_k] &= 0 \quad \text{for all } k.
\end{aligned}$$

Combining this with what is known about the marginal distribution of  $\mathbf{y}$  based on properties of normal distributions and letting  $\mathbf{t} = (1, \dots, n)$ , we have:

$$\mathbf{y} \sim \text{normal} \left( \mathbf{a} \circ \left( \phi^{t-1} \mu + \tilde{\mathbf{Z}} \mathbf{v} \right) + \mathbf{Z} \boldsymbol{\gamma}, \begin{pmatrix} \sigma_v^2 & \mathbf{0}' \\ \mathbf{0} & \sigma_w^2 \boldsymbol{\Omega} \circ (\mathbf{a} \mathbf{a}') + \sigma_v^2 \mathbf{I}_{n-1} \end{pmatrix} \right),$$

where ‘ $\circ$ ’ refers to the elementwise Hadamard product,  $\tilde{\mathbf{z}}_t = (\sum_{i=1}^t \phi^{t-i} \mathbf{z}'_i)$ , and  $\omega_{st} = \sigma_w^2 \sum_{i=2}^{\min\{s+1, t+1\}} \phi^{t+s+2(1-i)}$ . In the special case where  $a_t = a$ ,

$$\mathbf{y} \sim \text{normal} \left( \phi^{t-1} (a\mu) + \tilde{\mathbf{Z}} (a\mathbf{v}) + \mathbf{Z} \boldsymbol{\gamma}, \begin{pmatrix} \sigma_v^2 & \mathbf{0}' \\ \mathbf{0} & a^2 \sigma_w^2 \boldsymbol{\Omega} + \sigma_v^2 \mathbf{I}_{n-1} \end{pmatrix} \right).$$

Note that when  $\phi = 0$ ,  $\tilde{\mathbf{Z}} = \mathbf{Z}$ . This model assumes that the observations  $y_t$  are centered around a rescaled multiple of the starting value  $\mu$ , a weighted sum of the past and present covariate values, and the present covariate value. Note that  $a$ ,  $\phi$ ,  $\sigma_w^2$ , and  $\sigma_v^2$  are subject to the constraints that  $a^2 \sigma_w^2 \boldsymbol{\Omega} + \sigma_v^2 \mathbf{I}_{n-1}$  and  $\boldsymbol{\Omega}$  are positive definite, to ensure that the marginal likelihood and conditional distributions of the states are well defined.