

# Time Series and Sequence Learning

## Lecture 6 – Learning of State Space Models

---

Johan Alenlöv, Linköping University

2021-09-14

## Summary of Lecture 5: Trend component

A  $k - 1$ th order polynomial trend model  $\Delta^k \mu_t = \zeta_t$  can be written as

$$\alpha_t = \begin{bmatrix} c_1 & c_2 & \cdots & c_{k-1} & c_k \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix} \alpha_{t-1} + \begin{bmatrix} 1 \\ 0 \\ 0 \\ \cdots \\ 0 \end{bmatrix} \zeta_t,$$
$$\mu_t = \begin{bmatrix} 1 & 0 & 0 & \vdots & 0 \end{bmatrix} \alpha_t,$$

where the **state vector** is

$$\alpha_t = \begin{bmatrix} \mu_t & \mu_{t-1} & \cdots & \mu_{t-k+1} \end{bmatrix}^T$$

and  $c_i = (-1)^{i+1} \binom{k}{i}$ .

## Summary of Lecture 5: Seasonal component

A  $s$  period seasonal model,  $\sum_{j=0}^{s-1} \gamma_{t-j} = \omega_t$ , can be written a

$$\alpha_t = \begin{bmatrix} -1 & -1 & \cdots & -1 & -1 \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix} \alpha_{t-1} + \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \omega_t$$
$$\gamma_t = \begin{bmatrix} 1 & 0 & 0 & \vdots & 0 \end{bmatrix} \alpha_t,$$

where the **state vector** is

$$\alpha_t = \begin{bmatrix} \gamma_t & \gamma_{t-1} & \cdots & \gamma_{t-s+2} \end{bmatrix}^T.$$

# Summary of Lecture 5: Structural time series

A general structural time series model

$$y_t = \mu_t + \gamma_t + \varepsilon_t$$

can be written in state space form using **block matrices**.

**State vector:**

$$\alpha_t = \begin{bmatrix} \mu_t & \mu_{t-1} & \cdots & \mu_{t-k+1} & \gamma_t & \gamma_{t-1} & \cdots & \gamma_{t-s+2} \end{bmatrix}^T$$

**State space model:**

$$\begin{aligned} \alpha_t &= \begin{bmatrix} T_{[\mu]} & \\ & T_{[\gamma]} \end{bmatrix} \alpha_{t-1} + \begin{bmatrix} R_{[\mu]} & \\ & R_{[\gamma]} \end{bmatrix} \eta_t, & \eta_t &\sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_\zeta^2 & 0 \\ 0 & \sigma_\omega^2 \end{bmatrix} \right), \\ y_t &= \begin{bmatrix} Z_{[\mu]} & Z_{[\gamma]} \end{bmatrix} \alpha_t + \varepsilon_t, & \varepsilon_t &\sim \mathcal{N}(0, \sigma_\varepsilon^2). \end{aligned}$$

# Summary of Lecture 5: The Kalman filter

For any  $s, t$ , denote by  $\hat{\alpha}_{t|s} = \mathbb{E}[\alpha_t | y_{1:s}]$  and  $P_{t|s} = \text{Cov}(\alpha_t | y_{1:s})$ .

**Thm.** For an LGSS model,  $p(\alpha_t | y_{1:s}) = \mathcal{N}(\alpha_t | \hat{\alpha}_{t|s}, P_{t|s})$ .

Of particular interest are:

- Filtering distribution,

$$p(\alpha_t | y_{1:t}) = \mathcal{N}(\alpha_t | \hat{\alpha}_{t|t}, P_{t|t}).$$

- (1-step) Predictive distributions,

$$p(\alpha_t | y_{1:t-1}) = \mathcal{N}(\alpha_t | \hat{\alpha}_{t|t-1}, P_{t|t-1}),$$

$$p(y_t | y_{1:t-1}) = \mathcal{N}(y_t | \hat{y}_{t|t-1}, F_{t|t-1}).$$

# Summary of Lecture 5: ARMA model in state space form

**State space formulation of ARMA:** Consider the ARMA( $p, q$ ) model,

$$y_t = \sum_{j=1}^p a_j y_{t-j} + \sum_{j=1}^q b_j \eta_{t-j} + \eta_t.$$

Let  $d = \max(p, q + 1)$  and define  $a_j = 0$  for  $j > p$  and  $b_j = 0$  for  $j > q$ . Then, an equivalent **state space form** is given by

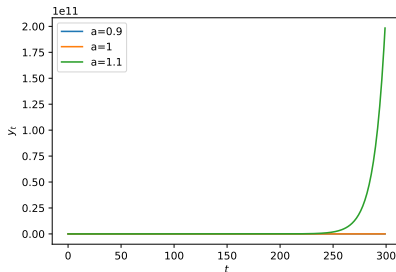
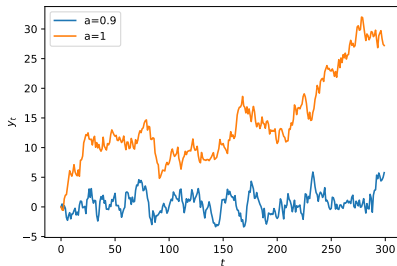
$$\alpha_t = \begin{bmatrix} a_1 & a_2 & \cdots & a_{d-1} & a_d \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix} \alpha_{t-1} + \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \eta_t,$$
$$y_t = \begin{bmatrix} 1 & b_1 & \cdots & b_{d-2} & b_{d-1} \end{bmatrix} \alpha_t$$

# Stability of state space models

---

## ex) Simulation of AR(1)

Simulation of  $y_t = ay_{t-1} + \varepsilon_t$



The AR(1) model is:

- Stable if  $|a| < 1 \Rightarrow$  converges to stationary
- Marginally stable if  $|a| = 1 \Rightarrow$  linear drift
- Unstable if  $|a| > 1 \Rightarrow$  exponential explosion



*Can this be generalized to an LGSS model?*

**AR(1):**

$$y_t = ay_{t-1} + \varepsilon_t$$

**LGSS:**

$$\alpha_t = T\alpha_{t-1} + R\eta_t$$

Intuitively, the state process is unstable if “size( $T$ )”  $> 1$ .

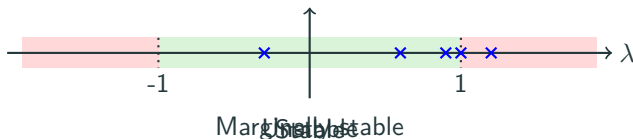
# Stability of state space model

**Thm.** A state space model of dimension  $d = \dim(\alpha_t)$  is:

- Stable iff  $|\lambda_j| < 1$ ,  $j = 1, \dots, d$ ,
- Marginally stable iff  $|\lambda_j| \leq 1$ ,  $j = 1, \dots, d$ ,
- Unstable iff  $|\lambda_j| > 1$  for any  $j = 1, \dots, d$ ,

where  $\lambda_j, j = 1, \dots, d$  are the **eigenvalues of  $T$** .

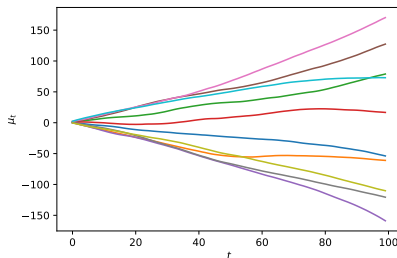
If all eigenvalues are real numbers,



## ex) Eigenvalues of linear trend model

Linear trend model:

$$\alpha_t = \overbrace{\begin{bmatrix} 2 & -1 \\ 1 & 0 \end{bmatrix}}^F \alpha_{t-1} + \overbrace{\begin{bmatrix} 1 \\ 0 \end{bmatrix}}^R \zeta_t$$
$$\mu_t = \underbrace{\begin{bmatrix} 1 & 0 \end{bmatrix}}_Z \alpha_t$$



Check for stability by computing the eigenvalues

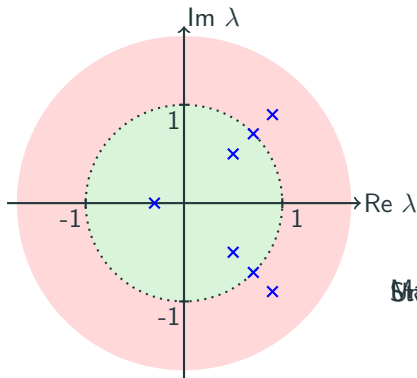
$$\text{eig}(F) \Rightarrow \lambda_1 = \lambda_2 = 1.$$

The trend model is **marginally stable!**

# Complex eigenvalues

**Note.** The eigenvalues can be complex numbers in general. Thus, the stability condition reads

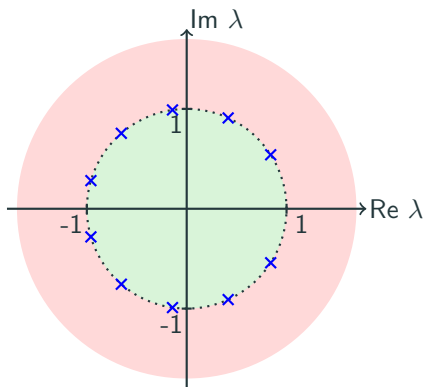
*all eigenvalues of  $T$  are within the unit circle in the complex plane.*



## ex) Eigenvalues of seasonal model

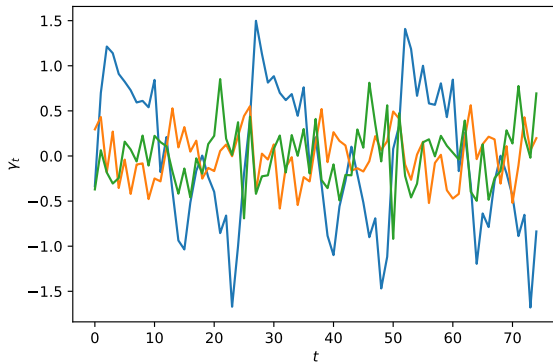
Seasonal model:

$$T = \begin{bmatrix} -1 & -1 & \cdots & -1 & -1 \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix}$$



The seasonal model is **marginally stable!**

## ex) Sample trajectories for seasonal model



The structural time series models that we have proposed are  
**designed to be marginally stable!**

Marginal stability results in desirable properties:

- Real eigenvalues  $\lambda_j = 1 \Rightarrow$  polynomial drift/trend.
- Complex eigenvalues with  $|\lambda_j| = 1 \Rightarrow$  periodicity/seasonality.

# Likelihood estimation

---



## Recap: The log-likelihood of the local level model

**The Local Level Model** given by

$$\begin{aligned}y_t &= \mu_t + \varepsilon_t, & \varepsilon_t &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_\varepsilon^2) \\ \mu_{t+1} &= \mu_t + \eta_{t+1}, & \eta_{t+1} &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_\eta^2)\end{aligned}$$

and initial distribution  $\mu_1 \sim \mathcal{N}(a_1, P_1)$

**Log-Likelihood** Given a sequence  $y_{1:n}$  the **log-likelihood** is given by

$$\begin{aligned}\ell(\theta) &= \log p_\theta(y_1) + \sum_{t=2}^n \log p_\theta(y_t | y_{1:t-1}) \\ &= -\frac{n}{2} \log 2\pi - \frac{1}{2} \sum_{t=1}^n \left( \log F_{t|t-1}(\theta) + \frac{(y_t - \hat{\mu}_{t|t-1}(\theta))^2}{F_{t|t-1}(\theta)} \right)\end{aligned}$$

# Log-likelihood of general SSM

**Def.** A **Linear Gaussian State-Space (LGSS)** model is given by:

$$\alpha_t = T\alpha_{t-1} + R\underline{\eta_t},$$

$$y_t = Z\alpha_t + \varepsilon_t$$

$$\eta_t \sim \mathcal{N}(0, \underline{Q}),$$
$$\varepsilon_t \sim \mathcal{N}(0, \underline{\sigma_\epsilon^2}),$$

and initial distribution  $\alpha_1 \sim \mathcal{N}(a_1, P_1)$ .

As in the local level model we write,

$$\ell(\theta) = \log p_\theta(y_1) + \sum_{t=2}^n \log p_\theta(y_t | y_{1:t-1})$$

Remains to find the distribution  $y_t | y_{1:t-1}$ .

## Calculating the log-likelihood

We again need to look at the distribution of  $y_t | y_{1:t-1}$ .

Gaussian distribution  $\Rightarrow$  find mean and variance.

$$\begin{aligned}\mathbb{E}[y_t | y_{1:t-1}] &= \mathbb{E}[Z\alpha_t + \varepsilon_t | y_{1:t-1}] = Z\hat{\alpha}_{t|t-1} = \hat{y}_{t|t-1} \\ \text{Var}[y_t | y_{1:t-1}] &= \text{Var}[Z\alpha_t + \varepsilon_t | y_{1:t-1}] \\ &= Z\text{Var}[\alpha_t | y_{1:t-1}]Z^T + \sigma_\epsilon^2 \\ &= ZP_{t|t-1}Z^T + \sigma_\epsilon^2 = F_{t|t-1}\end{aligned}$$

This gives us

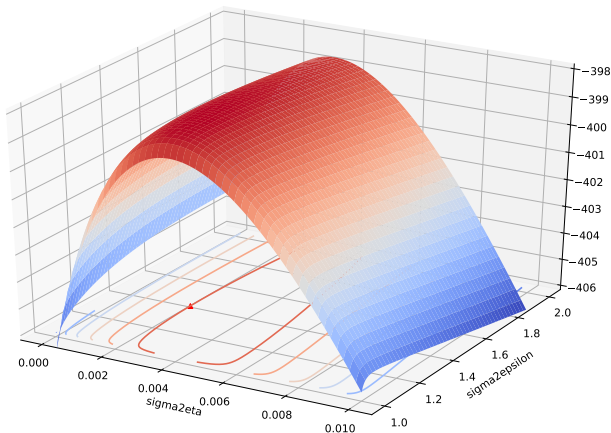
$$\log p_{\theta}(y_t | y_{1:t-1}) = \text{const} - \frac{1}{2} \left( \log |F_{t|t-1}| + (y_t - \hat{y}_{t|t-1})^T F_{t|t-1}^{-1} (y_t - \hat{y}_{t|t-1}) \right)$$

and

$$\ell(\theta) = \text{const} - \frac{1}{2} \sum_{t=1}^n \left( \log |F_{t|t-1}| + (y_t - \hat{y}_{t|t-1})^T F_{t|t-1}^{-1} (y_t - \hat{y}_{t|t-1}) \right)$$

## ex) Maximum likelihood in local level model

We return to the example of lecture 4.



Maximum found for  $\sigma_\epsilon^2 = 1.37$  and  $\sigma_\eta^2 = 0.002$ .

# Expectation-Maximization

---

## Another approach to parameter estimation

- Calculating the derivatives of  $\ell(\theta)$  is a hard problem.
- If we had access to  $\alpha_{1:n}$  it would be much easier.

$$\begin{aligned}\log p_{\theta}(\alpha_{1:n}, y_{1:n}) = \text{const.} - \frac{1}{2} \sum_{i=1}^n [\log |\sigma_{\epsilon}^2| + \log |Q| \\ + \underbrace{(y_i - Z\alpha_i)^T}_{\text{red underline}} \sigma_{\epsilon}^{-2} \underbrace{(y_i - Z\alpha_i)}_{\text{red underline}} \\ + \underbrace{(\alpha_i - T\alpha_{i-1})^T}_{\text{red underline}} R Q^{-1} R^T \underbrace{(\alpha_i - T\alpha_{i-1})}_{\text{red underline}}]\end{aligned}$$

$$\begin{aligned}\log p_{\theta}(\alpha_{1:n}, y_{1:n}) = \text{const.} - \frac{1}{2} \sum_{i=1}^n [\log |\sigma_{\epsilon}^2| + \log |Q| \\ + \underbrace{\epsilon_i^T}_{\text{red triangle}} \sigma_{\epsilon}^{-2} \underbrace{\epsilon_i}_{\text{red triangle}} + \underbrace{\eta_i^T}_{\text{red triangle}} Q^{-1} \underbrace{\eta_i}_{\text{red triangle}}] \\ \epsilon_i = y_i - Z\alpha_i \\ \eta_i = R^T(\alpha_i - T\alpha_{i-1})\end{aligned}$$

- Easy to take derivatives of this and maximize.
- Unfortunately we don't know  $\alpha_{1:t}$ , so can't use this directly.

# Expectation-Maximization

In the **Expectation Maximization** (EM) algorithm we alternate two steps,

1. E-step: Calculate  $Q(\theta, \tilde{\theta}) = \mathbb{E}[\log p_{\theta}(\alpha_{1:n}, y_{1:n}) \mid y_{1:n}, \tilde{\theta}]$
2. M-step: Find  $\theta^*$  that maximizes  $Q(\theta, \tilde{\theta})$ .

We have that, (using  $\mathbb{E}[x^T A x] = \text{tr}[A \Sigma] + m^T A m$  when  $x \sim \mathcal{N}(m, \Sigma)$ )

$$-\mathbb{E}[\log p_{\theta}(\alpha_{1:n}, y_{1:n}) \mid y_{1:n}, \tilde{\theta}] = \text{const.} - \frac{1}{2} \sum_{t=1}^n [\log |\sigma_{\epsilon}^2| + \log |Q| + \{\hat{\epsilon}_{t|n}^2 + \text{Var}[\epsilon_t \mid y_{1:n}]\sigma_{\epsilon}^{-2} + \text{tr}[\{\hat{\eta}_{t|n}\hat{\eta}_{t|n}^T + \text{Var}[\eta_t \mid y_{1:n}]\}Q^{-1}]],$$

where  $\hat{\epsilon}_{t|n}$ ,  $\text{Var}[\epsilon_t \mid y_{1:n}]$ ,  $\hat{\eta}_{t|n}$ , and  $\text{Var}[\eta_t \mid y_{1:n}]$  are the **smoothed** mean and variances of  $\epsilon_t$  and  $\eta_t$ .

To find  $\theta^*$  maximize  $Q(\theta, \tilde{\theta})$  by taking the derivative and set the derivative to zero.

## One Slide on the Proof

---



# The Smoothing Distribution

---

# The smoothing distribution

- **State smoothing** refers to the problem of estimating  $\alpha_t | y_{1:n}$  for  $t < n$ .
- Often separated into three classes:
  - **Fixed-interval smoothing**, when  $n$  is fixed.
  - **Fixed-point smoothing**, when  $t$  is fixed and  $n = t + 1, t + 2, \dots$
  - **Fixed-lag smoothing**, when  $t = n - \ell$ .
- In our case we are interested in the distributions

$$\underline{\eta_t | y_{1:n}} \quad \underline{\varepsilon_t | y_{1:n}},$$

this is known as **disturbance smoothing**.

- In the LGSS model the distributions will be **Gaussian**.

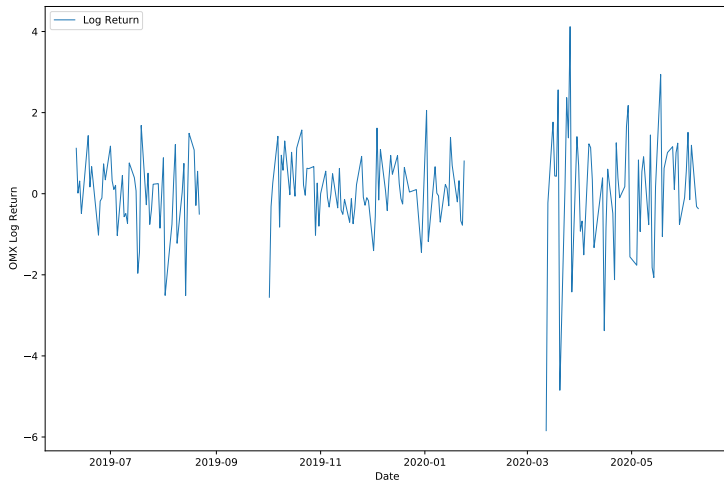
# State smoothing

- Smoothing is typically a **two-step** algorithm.
  - A **filter** is run in the forward direction ( $t = 1, 2, \dots, n$ )
  - A **smoother** is run in the backward direction ( $t = n, n - 1, \dots, 1$ )
- During the backward pass we will "correct" the filter distributions to the smoothing distributions.
- For  $\hat{\alpha}_{t|n}$  we get,

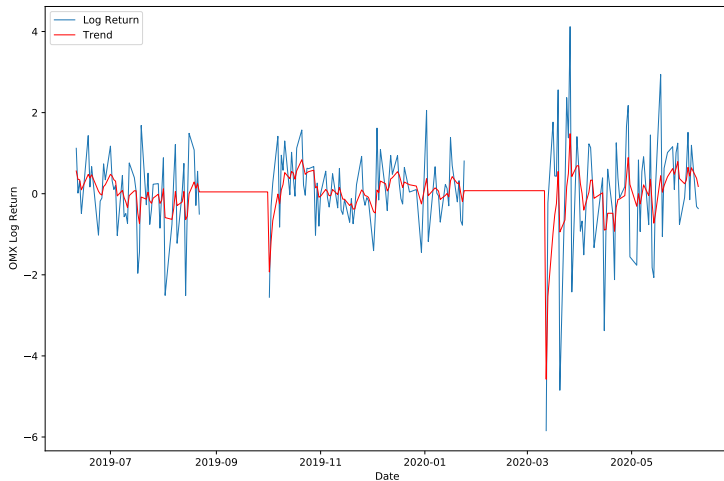
$$\begin{aligned}L_t &= T - TK_tZ \\r_{t-1} &= Z^T F_{t|t-1}^{-1}(y_t - \hat{y}_{t|t-1}) + L_t^T r_t \\ \hat{\alpha}_{t|n} &= \hat{\alpha}_{t|t-1} + P_{t|t-1} r_t.\end{aligned}$$

Initialized using  $r_n = 0$

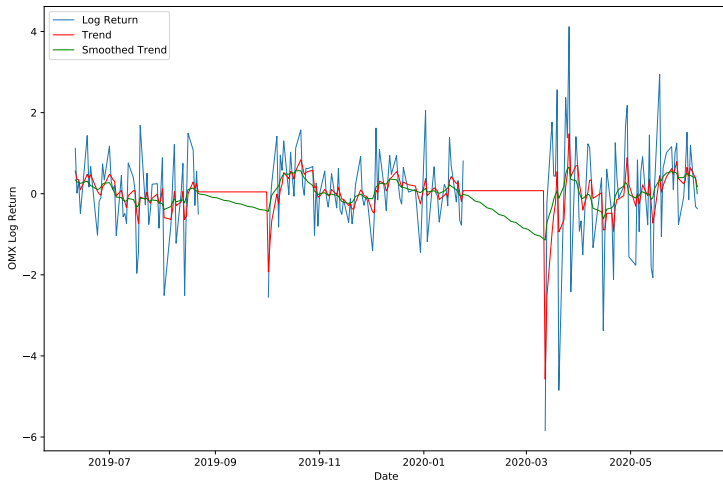
## ex) Filter vs. Smoothing distribution



## ex) Filter vs. Smoothing distribution



## ex) Filter vs. Smoothing distribution



---

## Kalman filter:

---

1. **Initialize:** Set  $\hat{\alpha}_{1|0} = a_1$  and  $P_{1|0} = P_1$ .

2. **for**  $t = 1, 2, \dots$

(a) **Measurement update:** // Skip if  $y_t$  is unavailable

$$\begin{aligned} \cdot \text{Predict } y_t: & \quad \begin{cases} \hat{y}_{t|t-1} = Z\hat{\alpha}_{t|t-1}, \\ F_{t|t-1} = ZP_{t|t-1}Z^T + \sigma_\varepsilon^2 \end{cases} \\ \cdot \text{Kalman gain:} & \quad K_t = P_{t|t-1}Z^T F_{t|t-1}^{-1} \\ \cdot \text{Update filter:} & \quad \begin{cases} \hat{\alpha}_{t|t} = \hat{\alpha}_{t|t-1} + K_t(y_t - \hat{y}_{t|t-1}), \\ P_{t|t} = (I - K_tZ)P_{t|t-1} \end{cases} \end{aligned}$$

(b) **Measurement update:**

$$\cdot \text{Predict } \alpha_t: \quad \begin{cases} \hat{\alpha}_{t+1|t} = T\hat{\alpha}_{t|t}, \\ P_{t+1|t} = TP_{t|t}T^T + \underline{RQR^T} \end{cases}$$

# State smoothing for general SSM

---

## State smoother:

---

1. **Initialize:** Run the **Kalman Filter** and save the Kalman gains and the predictive distributions.
2. **Initialize:** Set  $r_n = 0$  and  $N_n = 0$
3. **for**  $t = n, n - 1, \dots, 1$

$$\begin{aligned} \cdot \text{ Calculate:} & \quad \begin{cases} L_t = T - TK_tZ \\ r_{t-1} = Z^T F_{t|t-1}^{-1} (y_t - \hat{y}_{t|t-1}) + L_t^T r_t \\ N_{t-1} = Z^T F_{t|t-1}^{-1} Z + L_t^T N_t L_t \end{cases} \\ \cdot \text{ State smoothing:} & \quad \begin{cases} \hat{\alpha}_{t|n} = \hat{\alpha}_{t|t-1} + P_{t|t-1} r_{t-1} \\ P_{t|n} = P_{t|t-1} - P_{t|t-1} N_{t-1} P_{t|t-1} \end{cases} \end{aligned}$$

---



# Disturbance smoothing for general SSM

---

## Disturbance smoother:

---

1. **Initialize:** Run the **Kalman Filter** and save the Kalman gains and the predictive distributions.
2. **Initialize:** Set  $r_n = 0$  and  $N_n = 0$
3. **for**  $t = n, n - 1, \dots, 1$

· Calculate:

$$\begin{cases} C_t = T^T N_t T \\ D_t = F_{t|t-1}^{-1} + K_t^T C_t K_t \\ u_t = F_{t|t-1}^{-1} (y_t - \hat{y}_{t|t-1}) - K_t^T T^T r_t \end{cases}$$

· Observation noise:

$$\begin{cases} \hat{\varepsilon}_{t|n} = \sigma_\epsilon^2 u_t \\ \text{Var}[\varepsilon_t | y_{1:n}] = \sigma_\epsilon^2 - \sigma_\epsilon^2 D_t \sigma_\epsilon^2 \end{cases}$$

· State noise:

$$\begin{cases} \hat{\eta}_{t|n} = Q R^T r_t \\ \text{Var}[\eta_t | y_{1:n}] = Q - Q R^T N_t R Q \end{cases}$$

· Time update:

$$\begin{cases} r_{t-1} = Z^T u_t + T^T r_t \\ N_{t-1} = Z^T D_t Z + C_t - Z^T K_t^T C_t - C_t K_t Z \end{cases}$$

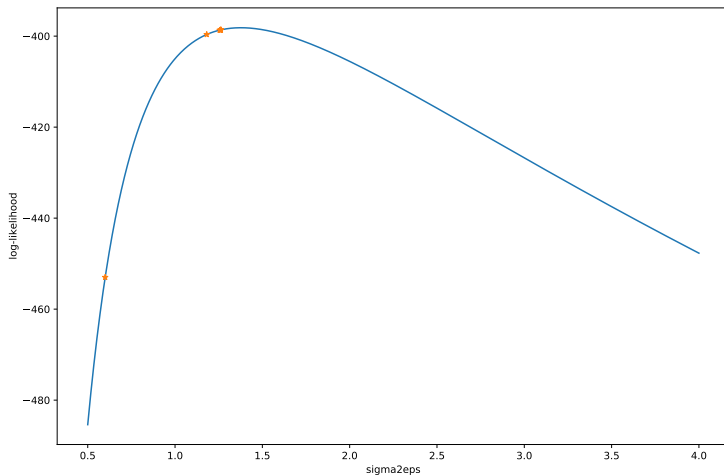
# The EM algorithm

- Set an initial parameter value  $\theta_0$ .
- For  $k = 0, 1, \dots$  do:
  1. Calculate the smoothing distribution using the **disturbance smoother** with the current parameter value  $\theta_k$ .
  2. Set  $\theta_{k+1} = \arg \max \mathcal{Q}(\theta, \theta_k)$ .

Until convergence.

## ex) The EM algorithm

Let's look at  $\sigma_\varepsilon^2$  for the local level model.



## A few concepts to summarize lecture 6:

- **Stability:** Results in a stationary state process (possibly after an initial transient). Corresponds to all eigenvalues of  $T$  being strictly within the unit circle.
- **Marginal stability:** Results in a state process that grows polynomially and/or shows a non-diminishing periodic pattern. Corresponds to some eigenvalues of  $T$  being *on* the unit circle.
- **Log-likelihood:** The log-likelihood for a LGSS can be calculated using the Kalman filter.
- **Expectation-Maximization:** Algorithm for maximum likelihood estimation. Iterates two steps, the **E-step** and **M-step**.
- **State Smoothing:** When estimating the hidden state  $\alpha_t$  conditioned on data  $y_{1:n}$  for  $n > t$ .
- **Disturbance Smoothing:** Estimation of the *noise variables*  $\eta_t$  and  $\varepsilon_t$  conditioned on  $y_{1:n}$  for  $n > t$ .