

# Modelling Non-linear and Non-stationary Time Series

## Chapter 7: State Space Models Maximum likelihood estimation

Henrik Madsen

Lecture Notes

September 2016

# Contents

- Non-linear State Space Model
- The Extended Kalman Filter
- ML Estimators for State Space Models
- Estimation of some Doubly Stochastic models
- The EM algorithm
- State Dependent Model (SDM)
- Generalized State Space Models

# State Space Models

- Due to the (first order) *Markov Property* of the state vector the state space representation is very useful for describing *non-stationary* and *time varying* systems.
- It will be seen that even linear state space models can be used to describe some *non-linear* processes.
- State space models provides a flexible approach to *adaptive estimation*.
- A class of state space models (SDM) will be introduced to describe most of the non-linear time series models that has been introduced.

# The non-linear, time-varying, discrete time state space model

Consider the *non-linear*, discrete time system

$$\mathbf{X}_{t+1} = \mathbf{f}(t, \mathbf{X}_t, \mathbf{U}_t) + \mathbf{v}_t \quad (1)$$

$$\mathbf{Y}_t = \mathbf{h}(t, \mathbf{X}_t, \mathbf{U}_t) + \mathbf{e}_t \quad (2)$$

where  $\{\mathbf{e}\}$  and  $\{\mathbf{v}\}$  are sequences of independent random variables with

$$\mathbf{E}[\mathbf{v}\mathbf{v}^T] = \mathbf{R}_{1,t}$$

$$\mathbf{E}[\mathbf{e}\mathbf{e}^T] = \mathbf{R}_{2,t}$$

$$\mathbf{E}[\mathbf{v}\mathbf{e}^T] = \mathbf{0}$$

# The Extended Kalman Filter

Given the previously defined non-linear state space models the *Extended Kalman Filter (EKF)* estimate of  $\mathbf{X}_{t+1}$  given  $\mathbf{Y}_0, \mathbf{Y}_1, \dots, \mathbf{Y}_t$  is given by

$$\hat{\mathbf{X}}_{t+1|t} = \mathbf{f}(t, \hat{\mathbf{X}}_t) + \mathbf{K}_t[\mathbf{Y}_t - \mathbf{h}(t, \hat{\mathbf{X}}_{t|t-1})] \quad (3)$$

where

$$\begin{aligned} \mathbf{K}_t &= \mathbf{F}(t, \hat{\mathbf{X}}_{t|t-1})\mathbf{P}_t\mathbf{H}^T(t, \hat{\mathbf{X}}_{t|t-1})[\mathbf{H}(t, \hat{\mathbf{X}}_{t|t-1})\mathbf{P}_t\mathbf{H}^T(t, \hat{\mathbf{X}}_{t|t-1}) + \mathbf{R}_{2,t}]^{-1} \\ \mathbf{P}_{t+1} &= \mathbf{F}(t, \hat{\mathbf{X}}_{t|t-1})\mathbf{P}_t\mathbf{F}^T(t, \hat{\mathbf{X}}_{t|t-1}) + \mathbf{R}_{1,t} \\ &\quad - \mathbf{K}_t[\mathbf{R}_{2,t} + \mathbf{H}(t, \hat{\mathbf{X}}_{t|t-1})\mathbf{P}_t\mathbf{H}^T(t, \hat{\mathbf{X}}_{t|t-1})]\mathbf{K}_t^T \end{aligned}$$

and where the Jacobians of  $\mathbf{f}$  and  $\mathbf{h}$  are given by

$$\begin{aligned} \mathbf{F}(t, \hat{\mathbf{X}}) &= \frac{\partial}{\partial \mathbf{X}} \mathbf{f}(t, \mathbf{X})|_{\mathbf{X}=\hat{\mathbf{X}}} \\ \mathbf{H}(t, \hat{\mathbf{X}}) &= \frac{\partial}{\partial \mathbf{X}} \mathbf{h}(t, \mathbf{X})|_{\mathbf{X}=\hat{\mathbf{X}}} \end{aligned}$$

This is often called the *discrete-discrete Extended Kalman Filter (EKF)*

# Maximum likelihood estimation

- Consider the linear stochastic state space model:

$$\mathbf{X}_t = \mathbf{A}_t \mathbf{X}_{t-1} + \mathbf{B}_t \mathbf{u}_{t-1} + \mathbf{e}_{1,t}, \quad (4)$$

$$\mathbf{Y}_t = \mathbf{C}_t \mathbf{X}_t + \mathbf{e}_{2,t}, \quad (5)$$

where  $\{\mathbf{e}_{1,t}\}$  and  $\{\mathbf{e}_{2,t}\}$  are mutually uncorrelated normal distributed white noise sequences with variance  $\Sigma_{1,t}$  and  $\Sigma_{2,t}$ , respectively.

- Assume that  $\dim \mathbf{Y} = m$  and that the model is **global identifiable**.
- The parameters are embedded in the matrices  $\mathbf{A}_t$ ,  $\mathbf{B}_t$ ,  $\mathbf{C}_t$ ,  $\Sigma_{1,t}$  and  $\Sigma_{2,t}$ , which might be *time varying*.
- As an example a missing observation at time  $t$  can be handled by putting

$$\Sigma_{2,t} = \infty \quad (6)$$

# ML estimates of state space models

Let  $\theta$  denote the unknown parameters, and  $\mathcal{Y}_N = (\mathbf{Y}_N, \dots, \mathbf{Y}_0)$  all  $\mathcal{N}$  available observations in the time series of length  $N$  ( $\mathcal{N} < N$ ).

The *conditional likelihood function* (conditioned on  $\mathbf{Y}_0$ ) is

$$L(\theta; \mathcal{Y}_N) = f(\mathcal{Y}_N | \theta) \quad (7)$$

$$= f(\mathbf{Y}_N | \mathcal{Y}_{N-1}, \theta) f(\mathbf{Y}_{N-1} | \mathcal{Y}_{N-2}, \theta) \cdots f(\mathbf{Y}_1 | \mathbf{Y}_0, \theta). \quad (8)$$

# ML estimation for state space models

- Since the model is linear and all the random variables are Gaussian it follows that  $\mathbf{Y}_{t+1}|\mathcal{Y}_t$  is Gaussian too.
- This conditional Gaussian density is characterized by the conditional mean and the conditional variance.
- Let us introduce the *one-step prediction error* (the *innovation*) and its covariance:

$$\tilde{\mathbf{Y}}_{t+1|t} = \mathbf{Y}_{t+1} - \hat{\mathbf{Y}}_{t+1|t}, \quad (9)$$

$$\begin{aligned} \mathbf{R}_{t+1} &= \mathbf{V} [\tilde{\mathbf{Y}}_{t+1|t}] = \mathbf{V} [\mathbf{Y}_{t+1}|\mathcal{Y}_t] = \Sigma_{t+1|t}^{yy} \\ &= \mathbf{C}_t \Sigma_{t+1|t}^{xx} \mathbf{C}_t^T + \Sigma_{2,t}. \end{aligned} \quad (10)$$

- The conditional mean and covariance are simply calculated by an ordinary Kalman Filter.



# ML estimates

- Given the mean and variance of  $\mathbf{Y}_{t+1}|\mathcal{Y}_t$  we are able to state the conditional density function

$$f(\mathbf{Y}_{t+1}|\mathcal{Y}_t) = [(2\pi)^m \det \mathbf{R}_{t+1}]^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} \tilde{\mathbf{Y}}_{t+1|t}^T \mathbf{R}_{t+1}^{-1} \tilde{\mathbf{Y}}_{t+1|t} \right] \quad (11)$$

- Using this conditional density function the conditional likelihood function writes

$$L(\boldsymbol{\theta}; \mathcal{Y}_{\mathcal{N}}) = \prod_{i=1}^{\mathcal{N}} [(2\pi)^m \det \mathbf{R}_i]^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} \tilde{\mathbf{Y}}_{i|i-1}^T \mathbf{R}_i^{-1} \tilde{\mathbf{Y}}_{i|i-1} \right], \quad (12)$$

- As starting values we can take  $\hat{\mathbf{X}}_{0|0} = 0$  and  $\boldsymbol{\Sigma}_{0|0}^{xx} = \alpha \mathbf{I}$  where  $\alpha$  is 'large'. Alternatively we can estimate the starting values using the ML approach.

# ML estimates of state space models

- The maximization of the likelihood function is equivalent to maximization of

$$\log L(\boldsymbol{\theta}; \mathcal{Y}_{\mathcal{N}}) = -\frac{1}{2} \sum_{i=1}^{\mathcal{N}} [\log \det \mathbf{R}_i + \tilde{\mathbf{Y}}_{i|i-1}^T \mathbf{R}_i^{-1} \tilde{\mathbf{Y}}_{i|i-1}] + \text{const.} \quad (13)$$

- The ML-estimate of  $\boldsymbol{\theta}$  is the argument which maximizes this, i.e.

$$\hat{\boldsymbol{\theta}} = \arg \left\{ \max_{\boldsymbol{\theta}} \log L(\boldsymbol{\theta}; \mathcal{Y}_{\mathcal{N}}) \right\}. \quad (14)$$

- For the maximization a numerical method must be used.

# Uncertainty of ML estimates

- As an approximation of the covariance of the parameter estimates we can use

$$\mathbf{V} \left[ \hat{\boldsymbol{\theta}} \right] \simeq -\mathbf{H}^{-1}, \quad (15)$$

where

$$\{\mathbf{H}\}_{ij} = \left. \frac{\partial^2 \log L(\boldsymbol{\theta}; \mathcal{Y}_{\mathcal{N}})}{\partial \theta_i \partial \theta_j} \right|_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}}. \quad (16)$$

- Note that  $\mathbf{H}$  is used as an approximation of the Fisher information matrix.

## ML estimates under stationary conditions

It is straight forward to see that in a *stationary situation* the problem reduces to the following likelihood function

$$L(\boldsymbol{\theta}; \mathcal{Y}_N) = \prod_{t=1}^N [(2\pi)^m \det \boldsymbol{\Sigma}]^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} \tilde{\mathbf{Y}}_{t|t-1}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Y}}_{t|t-1} \right] \quad (17)$$

$$= [(2\pi)^m \det \boldsymbol{\Sigma}]^{N/2} \exp \left[ -\frac{1}{2} \sum_{t=1}^N \tilde{\mathbf{Y}}_{t|t-1}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Y}}_{t|t-1} \right] \quad (18)$$

The assumption about stationarity implies that the matrices  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ ,  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$  are constant. This also means that missing observations are not allowed.

## $\Sigma$ known

In this case maximization of the likelihood function is equivalent to minimization of  $S_1(\boldsymbol{\theta})$  where

$$S_1(\boldsymbol{\theta}) = \frac{1}{N} \sum_{t=1}^N \tilde{\mathbf{Y}}_{t|t-1}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Y}}_{t|t-1} \quad (19)$$

where  $\tilde{\mathbf{Y}}_{t+1|t} = \mathbf{Y}_{t+1} - \hat{\mathbf{Y}}_{t+1|t}$  is the one-step prediction error. Let us introduce the *sample covariance of the predictions errors*, i.e.

$$\mathbf{D}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{t=1}^N \tilde{\mathbf{Y}}_{t|t-1} \tilde{\mathbf{Y}}_{t|t-1}^T \quad (20)$$

then the *cost function*  $S_1$  can be written

$$S_1(\boldsymbol{\theta}) = \text{tr} \boldsymbol{\Sigma}^{-1} \mathbf{D}(\boldsymbol{\theta}) \quad (21)$$

## $\Sigma$ unknown

- In this case the maximization of the likelihood function above is equivalent to minimization of

$$S(\boldsymbol{\theta}, \boldsymbol{\Sigma}) = \frac{1}{2}N \log \det \boldsymbol{\Sigma} + \frac{1}{2} \sum_{t=1}^N \tilde{\mathbf{Y}}_{t|t-1}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Y}}_{t|t-1} \quad (22)$$

- Differentiating with respect to  $\boldsymbol{\Sigma}$  gives

$$\frac{\partial S}{\partial \boldsymbol{\Sigma}} = \frac{N}{2} \boldsymbol{\Sigma}^{-1} - \frac{1}{2} \boldsymbol{\Sigma}^{-1} \left( \sum_{t=1}^N \tilde{\mathbf{Y}}_{t|t-1} \tilde{\mathbf{Y}}_{t|t-1}^T \right) \boldsymbol{\Sigma}^{-1} \quad (23)$$

which equals zero (ie. ML estimator) for

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{t=1}^N \tilde{\mathbf{Y}}_{t|t-1} \tilde{\mathbf{Y}}_{t|t-1}^T = \mathbf{D}(\boldsymbol{\theta}) \quad (24)$$

i.e. for any given value of  $\boldsymbol{\theta}$  then (24) minimizes (22). The problem can then be reduced by the constraint  $\boldsymbol{\Sigma} = \mathbf{D}(\boldsymbol{\theta})$ .

## Σ unknown

By substituting this constraint into the above cost function gives

$$S_c(\boldsymbol{\theta}) = \frac{1}{2} \log \det \mathbf{D}(\boldsymbol{\theta}) + \frac{1}{2} \sum_{t=1}^N \tilde{\mathbf{Y}}_{t|t-1}^T \mathbf{D}^{-1}(\boldsymbol{\theta}) \tilde{\mathbf{Y}}_{t|t-1} \quad (25)$$

or

$$S_c(\boldsymbol{\theta}) = \frac{1}{2} \log \det \mathbf{D}(\boldsymbol{\theta}) + \frac{1}{2} \text{tr} \left[ \left( \sum_{t=1}^N \tilde{\mathbf{Y}}_{t|t-1} \tilde{\mathbf{Y}}_{t|t-1}^T \right) \mathbf{D}^{-1}(\boldsymbol{\theta}) \right] \quad (26)$$

$$= \frac{1}{2} \log \det \mathbf{D}(\boldsymbol{\theta}) + \frac{N}{2} \text{tr} I_m \quad (27)$$

$$= \frac{1}{2} \log \det \mathbf{D}(\boldsymbol{\theta}) + \frac{Nm}{2} \quad (28)$$

Since the last term is constant the ML estimate is found by minimizing the *cost function*

$$S_2(\boldsymbol{\theta}) = \log \det \mathbf{D}(\boldsymbol{\theta}) \quad (29)$$

# Estimation of some doubly stochastic models

- Consider the *non-linear* process:

$$Y_t = \Phi_t Y_{t-1} + \epsilon_t \quad (30)$$

$$\Phi_t - \mu = \phi(\Phi_{t-1} - \mu) + \zeta_t \quad (31)$$

where  $\{\epsilon_t\}$  and  $\{\zeta_t\}$  are mutually uncorrelated Gaussian white noise processes with variances  $\sigma_\epsilon^2$  and  $\sigma_\zeta^2$ , respectively.

- This model is some times referred to as an AR(1)-AR(1) model, since both the model for  $Y_t$  and the embedded model for the parameter variation is an AR(1) model.
- The parameters of the non-linear model are  $(\phi, \mu, \sigma_\epsilon^2, \sigma_\zeta^2)$ .



## Estimation of some doubly stochastic models

By choosing a different parameterization of the model it can, however, be written as a linear state space model. The model for the parameter variation can be written

$$\Phi_t = \phi\Phi_{t-1} + \mu(1 - \phi) + \zeta_t$$

Hence, by introducing  $\delta = \mu(1 - \phi)$  the variation of both  $\{Y_t\}$  and  $\{\Phi_t\}$  can be described by the linear state space model

$$\begin{pmatrix} \Phi_t \\ \delta_t \end{pmatrix} = \begin{pmatrix} \phi & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \Phi_{t-1} \\ \delta_{t-1} \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} \zeta_t \quad (32)$$

$$Y_t = (Y_{t-1}, 0) \begin{pmatrix} \Phi_t \\ \delta_t \end{pmatrix} + \epsilon_t \quad (33)$$

Note that the new parameter is introduced as a state variable, and that there is a unique relation between the new parameters  $(\phi, \delta, \sigma_\epsilon^2, \sigma_\zeta^2)$  and the parameters of the model (30).

## A generalization

The principle outlined above is readily generalized. It is, for instance, rather easy to see that the following doubly stochastic model (an AR(1)-AR(2) model)

$$Y_t = \Phi_t Y_{t-1} + \epsilon_t \quad (34)$$

$$\Phi_t - \mu = \phi_1(\Phi_{t-1} - \mu) + \phi_2(\Phi_{t-1} - \mu) + \zeta_t \quad (35)$$

can be written as the following linear state space model

$$\begin{pmatrix} \Phi_t \\ \Phi_{t-1} \\ \delta_t \end{pmatrix} = \begin{pmatrix} \phi_1 & \phi_2 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \Phi_{t-1} \\ \Phi_{t-2} \\ \delta_{t-1} \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \zeta_t \quad (36)$$

$$Y_t = (Y_{t-1}, 0, 0) \begin{pmatrix} \Phi_t \\ \Phi_{t-1} \\ \delta_t \end{pmatrix} + \epsilon_t \quad (37)$$

where  $\delta = \mu(1 - \phi_1 - \phi_2)$ .

## An example

The above outlined class of doubly stochastic processes provides a rather flexible description of non-linear phenomena.

Consider as an example a doubly stochastic process with two different parameter sets as described here:

Parameter	$\phi$	$\mu$	$\sigma_{\zeta}^2$	$\sigma_{\epsilon}^2$
True values No. 1	0.85	0.80	$0.10^2$	$4.0^2$
True values No. 2	0.95	0.80	$0.10^2$	$4.0^2$

# An example (cont.)

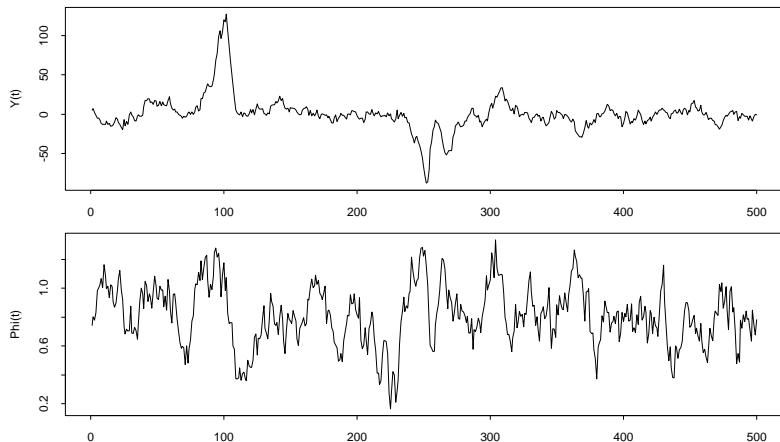


Figure : A simulation of a doubly stochastic process -  $\phi = 0.85$

# An example (cont.)

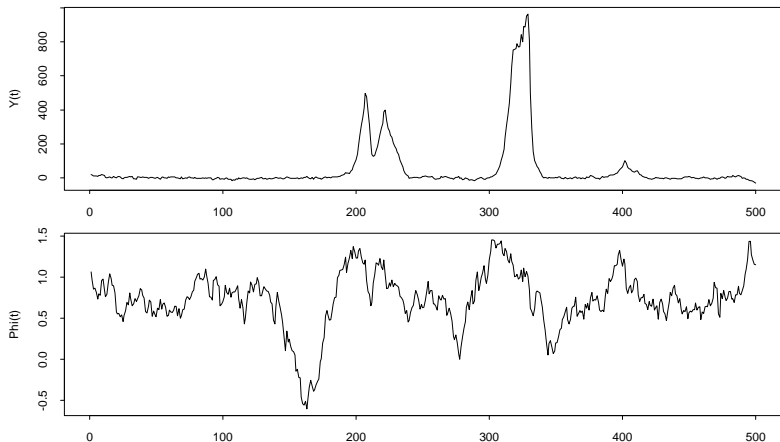


Figure : A simulation of a doubly stochastic process -  $\phi = 0.95$ .

# Extended Kalman Filter for parameter estimation

Idea: *Include the unknown parameters in the state vector* and then use a method for state estimation. The resulting state space model is in general a non-linear state space model, and the *extended Kalman filter* is then most often used for the state estimation.

Consider the following state space model

$$\begin{aligned} \mathbf{X}_{t+1} &= \mathbf{A}(\boldsymbol{\theta})\mathbf{X}_t + \mathbf{B}(\boldsymbol{\theta})\mathbf{u}_t + \mathbf{v}_t \\ \mathbf{Y}_t &= \mathbf{C}(\boldsymbol{\theta})\mathbf{X}_t + \mathbf{e}_t \end{aligned}$$

where the parameter vector  $\boldsymbol{\theta}$  is *unknown*.

# Extended Kalman Filter for parameter estimation

By including  $(\theta = \theta_t)$  in a new state vector

$$\mathbf{Z}_t = \begin{pmatrix} \mathbf{X}_t \\ \theta_t \end{pmatrix}$$

this leads to the following non-linear state space model

$$\begin{aligned} \mathbf{Z}_{t+1} &= f(\mathbf{Z}_t, \mathbf{u}_t) + \begin{pmatrix} \mathbf{v}_t \\ 0 \end{pmatrix} \\ \mathbf{Y}_t &= h(\mathbf{Z}_t) + \mathbf{e}_t \end{aligned}$$

where

$$\begin{aligned} f(\mathbf{Z}_t, \mathbf{u}_t) &= \begin{pmatrix} \mathbf{A}(\theta)\mathbf{X}_t + \mathbf{B}(\theta)\mathbf{u}_t \\ \theta_t \end{pmatrix} \\ h(\mathbf{Z}_t) &= \mathbf{C}(\theta)\mathbf{X}_t \end{aligned}$$

Now an EKF state estimation will provide a simultaneous estimation of the original state vector and the parameter vector  $\theta$

# The EM algorithm

For some problems the likelihood function for an extended (complete) data set is more easy to formulate than the likelihood function for the (incomplete) data at hand, which is assumed to be a subset of the complete data set.

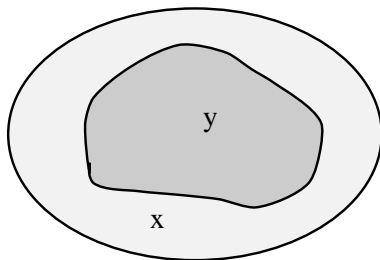


Figure : Complete (x) and incomplete data (y).

The incompleteness may be due to e.g. missing observations, or



The EM algorithm contains two steps:

- 1 The **E**xpectation step for estimating the complete likelihood function conditioning on the available data and the latest parameter estimates.
- 2 The **M**aximization step for finding the optimal parameter values in the current iteration.

Let  $k(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})$  be the conditional density of  $\mathbf{x}$  given  $\mathbf{y}$  and  $\boldsymbol{\theta}$ . Since  $\mathbf{y}$  is a subset of  $\mathbf{x}$  we have

$$k(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) = \frac{g(\mathbf{x}|\boldsymbol{\theta})}{f(\mathbf{y}|\boldsymbol{\theta})}.$$

or

$$\log f(\mathbf{y}|\boldsymbol{\theta}) = \log g(\mathbf{x}|\boldsymbol{\theta}) - \log k(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})$$

By multiplying both sides by the density  $p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(p)})$  and integrate with respect to  $\mathbf{x}$  we obtain the log likelihood

$$L(\boldsymbol{\theta}) = \log f(\mathbf{y}|\boldsymbol{\theta}) = Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)}) - H(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)})$$

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)}) = \text{E} \left[ \log g(\mathbf{x}|\boldsymbol{\theta}) | \mathbf{y}, \boldsymbol{\theta}^{(p)} \right]$$

$$H(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)}) = \text{E} \left[ \log k(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) | \mathbf{y}, \boldsymbol{\theta}^{(p)} \right].$$

# The EM-iteration

Now we can define the EM-iteration

- *E-step*: Calculate  $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)}) = \mathbb{E}_p[\log L(\mathbf{x}, \boldsymbol{\theta}|y_1, \dots, y_n, \boldsymbol{\theta}^{(p)})]$ , as a function of  $\boldsymbol{\theta}$ .
- *M-step*: Determine  $\boldsymbol{\theta}^{(p+1)}$  such that  $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)})$  is maximized.

That is the EM algorithm. The EM algorithm thus finds the maximum of the likelihood function iteratively.

## Sketch of the proof

Jensen inequality states that for a random variable  $X$  and a convex function  $\phi(x)$

$$\mathbb{E}[\phi(X)] \geq \phi(\mathbb{E}[X])$$

That is (here: concave function)

$$\begin{aligned} \mathbb{E} \left[ \log k(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) | \mathbf{y}, \boldsymbol{\theta}^{(p)} \right] &\leq \log k(\mathbf{x}^{(p)} | \mathbf{y}, \boldsymbol{\theta}^{(p)}) \\ H(\boldsymbol{\theta}, \boldsymbol{\theta}^{(p)}) &\leq H(\boldsymbol{\theta}^{(p)}, \boldsymbol{\theta}^{(p)}) \end{aligned}$$

And hence

$$L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^{(p)}) = Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(p)}) - Q(\boldsymbol{\theta}^{(p)}, \boldsymbol{\theta}^{(p)}) + H(\boldsymbol{\theta}^{(p)}, \boldsymbol{\theta}^{(p)}) - H(\boldsymbol{\theta}, \boldsymbol{\theta}^{(p)}) \geq 0$$

if  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(p)})$  is maximized.

Hence, the EM algorithm iteration never decreases the log likelihood.

*The key result is that in order to maximize  $L$  we only need to maximize  $Q$ .*

# EM-estimation of NLAR(1) model

$$Y_t = h(Y_{t-1}, \theta) + \epsilon_t \quad (38)$$

where  $\{\epsilon_t\}$  is i.i.d.  $N(0,1)$ .

Assume that the available observations are

$$\mathcal{Y}_N = (Y_1, Y_2, \dots, Y_{j-1}, Y_{j+1}, \dots, Y_N) \quad (39)$$

i.e.  $Y_j$  is missing.

The  $i$ 'th iteration of the EM-algorithm is

- **E-step**: Find  $E[Y_j|\mathcal{Y}_N, \theta^{(i)}] = Y_j^{(i)}$ ,  $E[Y_j^2|\mathcal{Y}_N, \theta^{(i)}] = (Y_j^2)^{(i)}$ ,  $E[h(Y_j)|\mathcal{Y}_N, \theta^{(i)}] = h(Y_j)^{(i)}$ ,  $E[h^2(Y_j)|\mathcal{Y}_N, \theta^{(i)}] = h^2(Y_j)^{(i)}$ .

In the evaluation it can be used that (for the NLAR(1)-model)

$$p(Y_j|\mathcal{Y}_N) = \frac{p(Y_j|Y_{j-1})p(Y_{j+1}|Y_j)}{p(Y_{j+1}|Y_{j-1})} \quad (40)$$

Note also that  $p(Y_{j+1}|Y_j) = \phi(Y_{j+1} - h(Y_j))$  where  $\phi$  is the density function of  $N(0,1)$ .

- *M-step* : Minimize with respect to  $\theta$  the 'sum of squares'

$$\sum_{t=1}^N Y_t^2 - 2 \sum_{t=1}^N Y_t h(Y_{t-1}) + \sum_{t=1}^N h^2(Y_{t-1}) \quad (41)$$

where we replace  $Y_j, Y_j^2, h(Y_j), h^2(Y_j)$  by the values found under the E-step. This gives  $\theta^{(i+1)}$ .

# State-dependent models (SDM)

- Previously we introduced the general non-linear model,

$$X_t = h'(\epsilon_t, \epsilon_{t-1}, \dots) \quad (42)$$

which clearly is “infinite dimensional”, since infinitely many past variables are used.

- The idea behind the state-dependent models is to reduce this to a “finite dimensional” form.
- In a parsimonious finite description it might be a good idea to consider previous  $X$  and  $\epsilon$ -values simultaneously. Hence consider:

$$X_t = h(X_{t-1}, \dots, X_{t-p}, \epsilon_{t-1}, \dots, \epsilon_{t-q}) + \epsilon_t \quad (43)$$



# State-dependent models

- Define the “state vector”

$$\tilde{\mathbf{x}}_t = (\epsilon_{t-q+1}, \dots, \epsilon_t, X_{t-p+1}, \dots, X_t)^T$$

Priestly (1980, 1988) uses the term “state vector” although the definition does not in general lead to a *minimal realization*<sup>1</sup>. It is for instance well-known that for an ARMA( $p, q$ ) model the minimal realization has a state vector of dimension  $\max(p, q + 1)$ , whereas the state vector  $\tilde{\mathbf{x}}_t$  has the dimension  $(q + p)$ .

- After a first-order Taylor expansion of  $h$  about some fixed time point we obtain the *State-Dependent Model (SDM)* of order  $(p, q)$ :

$$X_t + \sum_{i=1}^p \phi_i(\tilde{\mathbf{x}}_{t-1})X_{t-i} = \mu(\tilde{\mathbf{x}}_{t-1}) + \epsilon_t + \sum_{i=1}^q \psi_i(\tilde{\mathbf{x}}_{t-1})\epsilon_{t-i} \quad (44)$$

---

<sup>1</sup>Recall that a realization of a state space model (represented by the triplet of matrices  $(A, B, C)$ ) is minimal iff the corresponding system is completely controllable and completely observable

## Some nonlinear models on SDM form

- *ARMA models*:  $\mu$ ,  $\{\phi_i\}$  and  $\{\psi_i\}$  constants.
- *Bilinear models*:  $\mu$ ,  $\{\phi_i\}$  constants, and

$$\psi_i(\tilde{\mathbf{x}}_{t-1}) = c_i + \sum_{j=1}^m b_{ji} X_{t-j} \quad ; \quad i = 1, \dots, q \quad (45)$$

By selecting  $q$  properly and by putting some of  $c_i$  or  $b_{ji}$  equal to zero the general bilinear model is obtained.

- *Threshold models*:

$$\{\psi_i\} = 0 \quad (46)$$

$$\mu_t = a_0^{(J_t)} \quad (47)$$

$$\phi_i(\tilde{\mathbf{x}}_{t-1}) = -a_i^{(J_t)} \quad \text{if} \quad X_{t-d} \in R^{(J_t)} \quad (48)$$

- *Exponential AR models*: For  $\mu = 0$ ,  $\{\psi_i\} = 0$  and

$$\phi_i(\tilde{\mathbf{x}}_{t-1}) = -(\alpha_i + \beta_i e^{-\delta X_{t-1}^2}) \quad i = 1, \dots, p \quad (49)$$

then (44) reduces to the EXPAR model.

## A state space from of SDM models

The “state vector”

$$\tilde{\mathbf{x}}_t = (\epsilon_{t-q+1}, \dots, \epsilon_t, X_{t-p+1}, \dots, X_t)^T$$

together with the innovation,  $\epsilon_{t+1}$ , determine completely the evolution from time  $t$  to time  $t + 1$ .

A formal state space representation is

$$\tilde{\mathbf{x}}_{t+1} = \boldsymbol{\mu}(\tilde{\mathbf{x}}_t) + \mathbf{F}(\tilde{\mathbf{x}}_t)\tilde{\mathbf{x}}_t + \tilde{\boldsymbol{\epsilon}}_{t+1} \quad (50)$$

$$X_t = \mathbf{C}\tilde{\mathbf{x}}_t \quad (51)$$

where

$$\boldsymbol{\mu}(\tilde{\mathbf{x}}_{t-1}) = (0, \dots, 0 \dot{:} 0, \dots, \boldsymbol{\mu}(\tilde{\mathbf{x}}))^T \quad (52)$$

$$\tilde{\boldsymbol{\epsilon}}_t = \epsilon_t(0, \dots, 1 \dot{:} 0, \dots, 1)^T \quad (53)$$

$$\mathbf{C} = (0, \dots, 0, 1) \quad (54)$$

$$(55)$$

The F matrix is

$$\mathbf{F}(\tilde{\mathbf{x}}_t) = \begin{bmatrix}
 0 & 1 & 0 & \dots & 0 & \vdots & 0 & \dots & \dots & \dots & 0 \\
 0 & 0 & 1 & \dots & 0 & \vdots & 0 & \dots & \dots & \dots & 0 \\
 \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots & & & & \vdots \\
 0 & 0 & 0 & \dots & 1 & \vdots & 0 & \dots & \dots & \dots & 0 \\
 0 & 0 & 0 & \dots & 0 & \vdots & 0 & \dots & \dots & \dots & 0 \\
 \dots & \dots & \dots & \dots & \dots & \vdots & \dots & \dots & \dots & \dots & \dots \\
 0 & \dots & \dots & \dots & 0 & \vdots & 0 & 1 & 0 & \dots & 0 \\
 0 & \dots & \dots & \dots & 0 & \vdots & 0 & 0 & 1 & \dots & 0 \\
 \vdots & & & & \vdots & \vdots & \vdots & \vdots & \vdots & & \vdots \\
 0 & \dots & \dots & \dots & 0 & \vdots & 0 & 0 & 0 & \dots & 1 \\
 \psi_q & \psi_{q-1} & & \dots & \psi_1 & \vdots & -\phi_p & -\phi_{p-1} & & \dots & -\phi_1
 \end{bmatrix}$$

(56)