



Time Series Analysis

Henrik Madsen

`hm@imm.dtu.dk`

Informatics and Mathematical Modelling
Technical University of Denmark
DK-2800 Kgs. Lyngby



Outline of the lecture

Recursive and adaptive estimation:

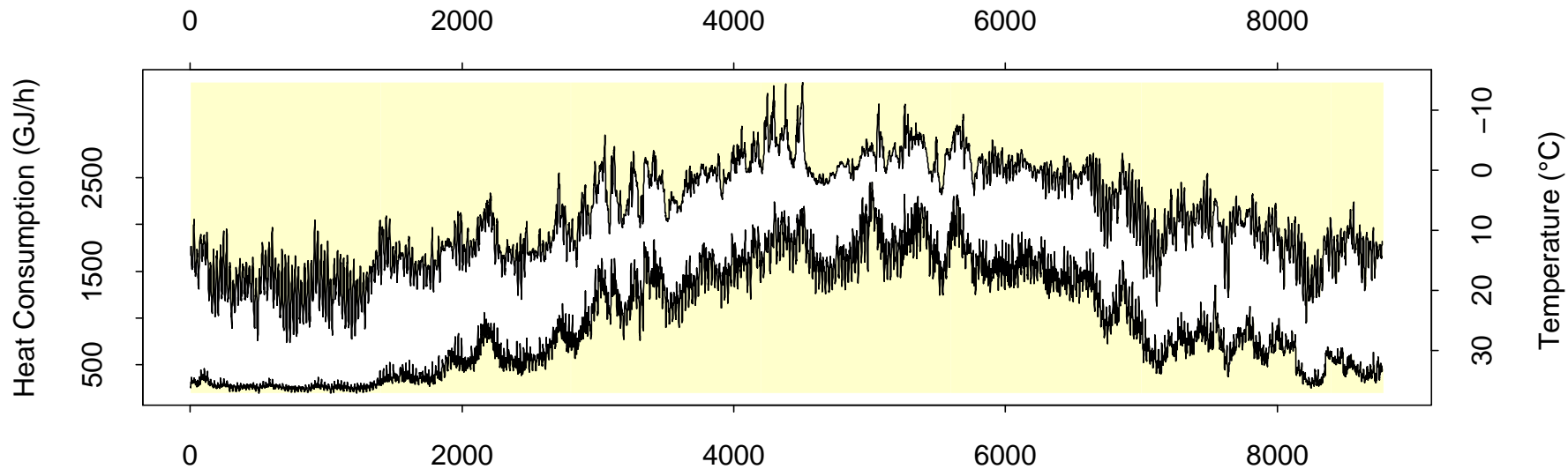
- Introduction to Chapter 11
- Recursive LS, Section 11.1

Cursory material:

- Model based adaptive estimation, Section 11.4



Why recursive and adaptive estimation?



- As time passes we get more information
- Models are approximations
- The best approximation change over time
- Makes it possible to produce software which learns as new data becomes available



Types of models considered

REG:

$$Y_t = \mu + \beta_1 U_{1,t} + \beta_2 U_{2,t} + \dots + \beta_m U_{m,t} + \varepsilon_t$$

FIR:

$$\begin{aligned} Y_t &= \mu + \omega(B)U_t + \varepsilon_t \\ &= \mu + \omega_0 U_t + \omega_1 U_{t-1} + \dots + \omega_s U_{t-s} + \varepsilon_t \end{aligned}$$

AR:

$$\begin{aligned} \phi(B)Y_t &= \mu + \varepsilon_t \Leftrightarrow \\ Y_t &= \mu - \phi_1 Y_{t-1} - \phi_2 Y_{t-2} - \dots - \phi_p Y_{t-p} + \varepsilon_t \end{aligned}$$

ARX:

$$\begin{aligned} \phi(B)Y_t &= \mu + \omega(B)U_t + \varepsilon_t \Leftrightarrow \\ Y_t &= \mu - \phi_1 Y_{t-1} - \dots - \phi_p Y_{t-p} + \omega_0 U_t + \dots + \omega_s U_{t-s} + \varepsilon_t \end{aligned}$$



Generic form of the models considered

$$\begin{aligned} Y_t &= \mathbf{x}_t^T \boldsymbol{\theta} + \varepsilon_t \\ &= \theta_1 x_{1,t} + \theta_2 x_{2,t} + \dots + \theta_\ell x_{\ell,t} + \varepsilon_t \end{aligned}$$

Example:

$$Y_t = \mu \cdot \underbrace{1}_{x_{1,t}} + \phi_2 \cdot \underbrace{(-Y_{t-2})}_{x_{2,t}} + \omega_1 \cdot \underbrace{U_{t-1}}_{x_{3,t}} + \varepsilon_t$$



LS-estimate at time t

Model:

$$Y_t = \mathbf{x}_t^T \boldsymbol{\theta} + \varepsilon_t$$

Data (\mathbf{x} may contain lagged values of the “real” input):

$$Y_1, Y_2, Y_3, Y_4, \dots, Y_{t-1}, Y_t$$
$$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \dots, \mathbf{x}_{t-1}, \mathbf{x}_t$$

LS-estimate based on t observations:

$$\hat{\boldsymbol{\theta}}_t = \arg \min_{\boldsymbol{\theta}} S_t(\boldsymbol{\theta})$$
$$S_t(\boldsymbol{\theta}) = \sum_{s=1}^t (Y_s - \mathbf{x}_s^T \boldsymbol{\theta})^2$$



LS-estimate at time t

Model:

$$Y_t = \mathbf{x}_t^T \boldsymbol{\theta} + \varepsilon_t$$

Data (\mathbf{x} may contain lagged values of the “real” input):

$$Y_1, Y_2, Y_3, Y_4, \dots, Y_{t-1}, Y_t$$
$$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \dots, \mathbf{x}_{t-1}, \mathbf{x}_t$$

LS-estimate based on t observations:

$$\hat{\boldsymbol{\theta}}_t = \arg \min_{\boldsymbol{\theta}} S_t(\boldsymbol{\theta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$
$$S_t(\boldsymbol{\theta}) = \sum_{s=1}^t (Y_s - \mathbf{x}_s^T \boldsymbol{\theta})^2$$



From one time step to the next (in an easy way)

The trick is to realize that:

$$\mathbf{R}_t = \mathbf{X}^T \mathbf{X} = \mathbf{x}_1 \mathbf{x}_1^T + \mathbf{x}_2 \mathbf{x}_2^T + \dots + \mathbf{x}_t \mathbf{x}_t^T = \sum_{s=1}^t \mathbf{x}_s \mathbf{x}_s^T$$

$$\mathbf{h}_t = \mathbf{X}^T \mathbf{Y} = \mathbf{x}_1 Y_1 + \mathbf{x}_2 Y_2 + \dots + \mathbf{x}_t Y_t = \sum_{s=1}^t \mathbf{x}_s Y_s$$

Where:

$$\mathbf{x}_t \mathbf{x}_t^T = \begin{bmatrix} x_{1,t} x_{1,t} & x_{1,t} x_{2,t} & \cdots & x_{1,t} x_{l,t} \\ x_{2,t} x_{1,t} & x_{2,t} x_{2,t} & \cdots & x_{2,t} x_{l,t} \\ \vdots & \vdots & \ddots & \vdots \\ x_{l,t} x_{1,t} & x_{l,t} x_{2,t} & \cdots & x_{l,t} x_{l,t} \end{bmatrix} \quad \mathbf{x}_t Y_t = \begin{bmatrix} x_{1,t} Y_t \\ x_{2,t} Y_t \\ \vdots \\ x_{l,t} Y_t \end{bmatrix}$$



From one time step to the next (cont'nd)

$$\hat{\theta}_t = \mathbf{R}_t^{-1} \mathbf{h}_t$$

$$\mathbf{R}_t = \sum_{s=1}^t \mathbf{x}_s \mathbf{x}_s^T = \mathbf{x}_t \mathbf{x}_t^T + \sum_{s=1}^{t-1} \mathbf{x}_s \mathbf{x}_s^T = \underline{\mathbf{x}_t \mathbf{x}_t^T + \mathbf{R}_{t-1}}$$

$$\mathbf{h}_t = \sum_{s=1}^t \mathbf{x}_s Y_s = \mathbf{x}_t Y_t + \sum_{s=1}^{t-1} \mathbf{x}_s Y_s = \underline{\mathbf{x}_t Y_t + \mathbf{h}_{t-1}}$$

Initialization:

- $\mathbf{R}_0 = \mathbf{0}$ (matrix of zeros)
- $\mathbf{h}_0 = \mathbf{0}$ (vector of zeros)
- Wait with $\hat{\theta}_t$ until \mathbf{R}_t is invertible



Other formulations

1. Eliminating h_t :

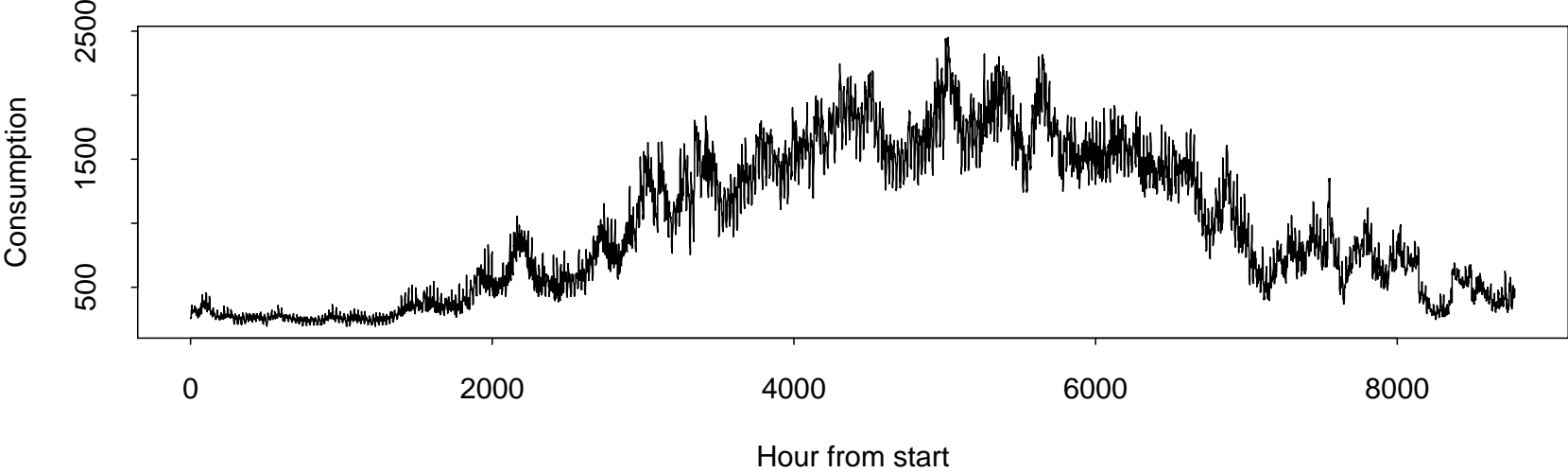
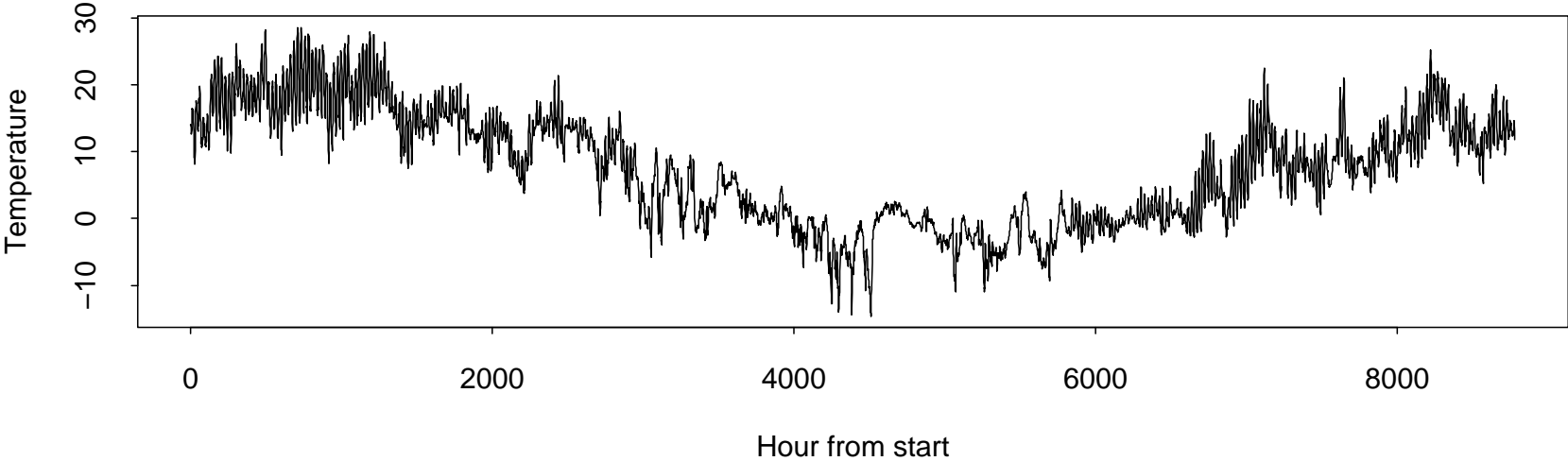
$$\begin{aligned} \mathbf{R}_t &= \mathbf{x}_t \mathbf{x}_t^T + \mathbf{R}_{t-1} \\ \hat{\boldsymbol{\theta}}_t &= \hat{\boldsymbol{\theta}}_{t-1} + \mathbf{R}_t^{-1} \mathbf{x}_t (Y_t - \mathbf{x}_t^T \hat{\boldsymbol{\theta}}_{t-1}) \end{aligned}$$

2. Eliminating h_t and avoiding matrix-inversion:

$$\begin{aligned} \mathbf{K}_t &= \frac{\mathbf{P}_{t-1} \mathbf{x}_t}{1 + \mathbf{x}_t^T \mathbf{P}_{t-1} \mathbf{x}_t} \\ \hat{\boldsymbol{\theta}}_t &= \hat{\boldsymbol{\theta}}_{t-1} + \mathbf{K}_t (Y_t - \mathbf{x}_t^T \hat{\boldsymbol{\theta}}_{t-1}) \\ \mathbf{P}_t &= \mathbf{P}_{t-1} - \frac{\mathbf{P}_{t-1} \mathbf{x}_t \mathbf{x}_t^T \mathbf{P}_{t-1}}{1 + \mathbf{x}_t^T \mathbf{P}_{t-1} \mathbf{x}_t} \end{aligned}$$

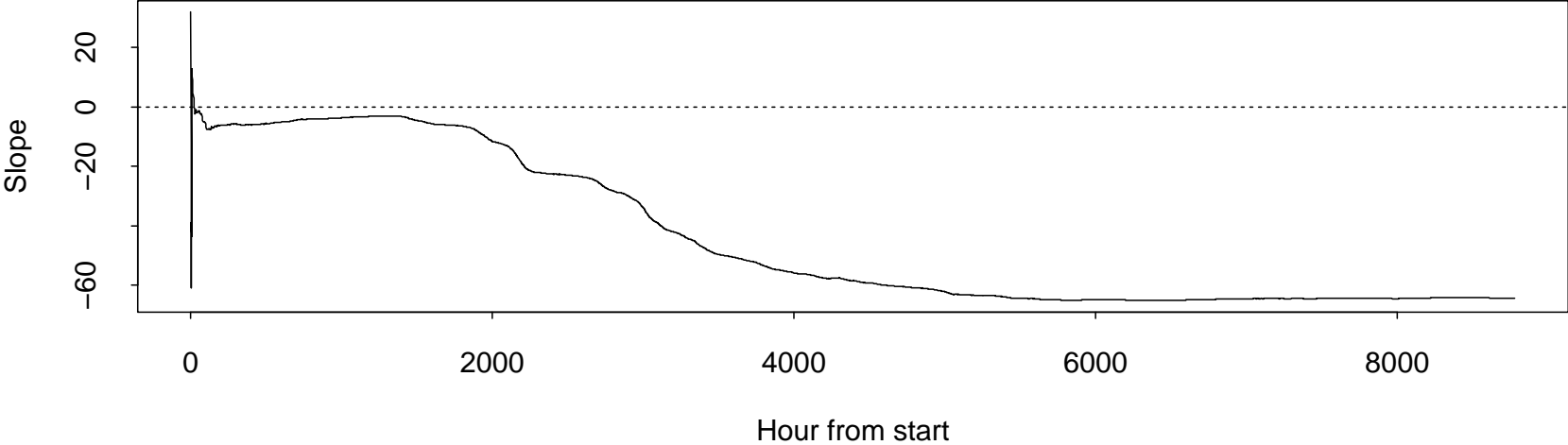
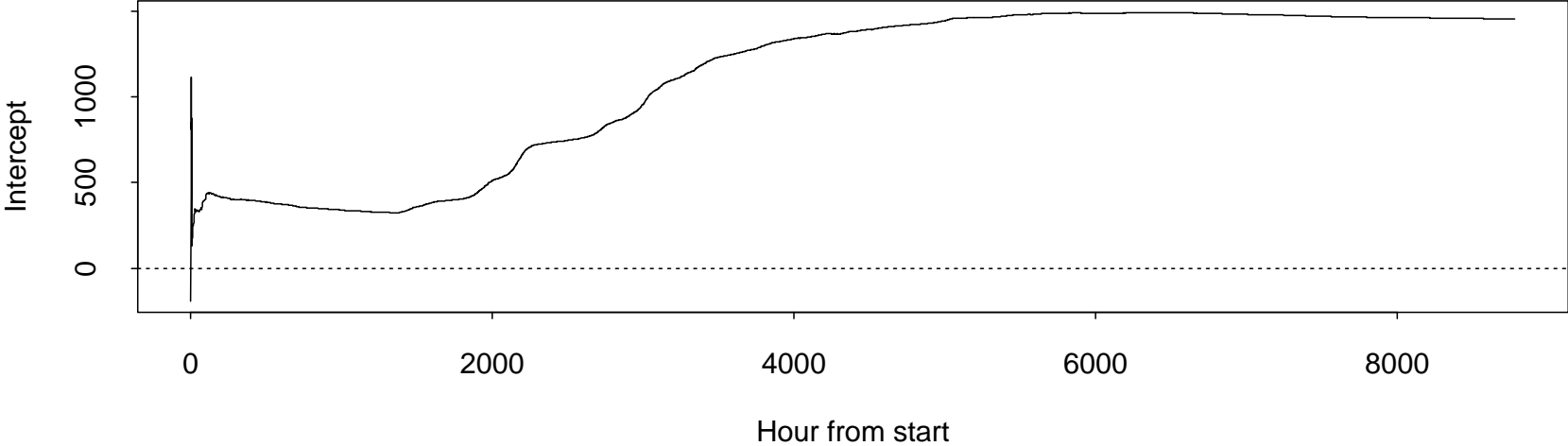


Example: $HC_t = \mu + \theta_1 T_t + \varepsilon_t$





Example: $HC_t = \mu + \theta_1 T_t + \varepsilon_t$





Forgetting old observations

- So far we have a way of updating the estimates as the data set grows
- If we want a method which forgets old observations we apply weights which start at 1 and goes to 0 when observations gets old:

$$\hat{\boldsymbol{\theta}}_t = \arg \min_{\boldsymbol{\theta}} S_t(\boldsymbol{\theta})$$
$$S_t(\boldsymbol{\theta}) = \sum_{s=1}^t \beta(t, s) (Y_s - \mathbf{x}_s^T \boldsymbol{\theta})^2$$

- $\beta(t, s)$ express how we assign weights to old observations



Forgetting old observations

- So far we have a way of updating the estimates as the data set grows
- If we want a method which forgets old observations we apply weights which start at 1 and goes to 0 when observations gets old:

$$\hat{\boldsymbol{\theta}}_t = \arg \min_{\boldsymbol{\theta}} S_t(\boldsymbol{\theta}) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{Y}$$

$$S_t(\boldsymbol{\theta}) = \sum_{s=1}^t \beta(t, s) (Y_s - \mathbf{x}_s^T \boldsymbol{\theta})^2$$

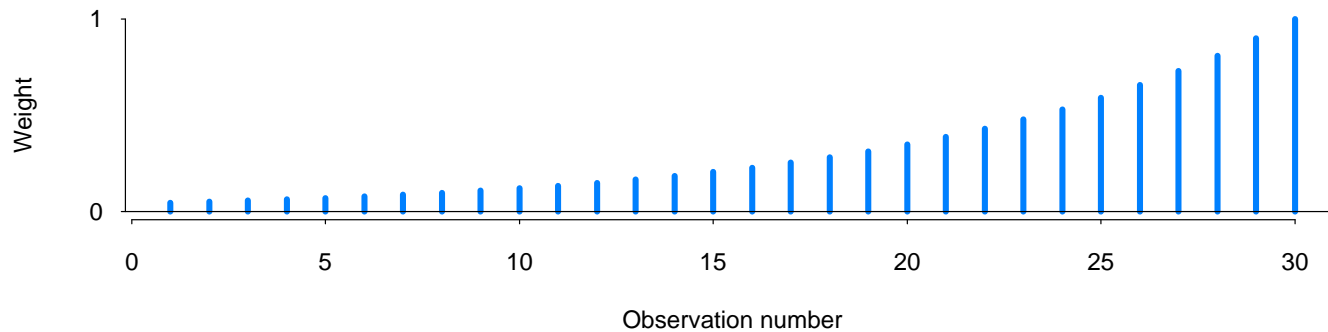
where $\mathbf{W} = \text{diag}(\beta(t, 1), \beta(t, 2), \dots, \beta(t, t - 1), 1)$

- $\beta(t, s)$ express how we assign weights to old observations



Exponential decay of weights

- Let's first consider $\beta(t, s) = \lambda^{t-s}$ ($0 < \lambda \leq 1$)
 - $\lambda = 1$: What we did with the previous algorithms
 - $\lambda < 1$: We forget in an exponential manner



- In the general case it turns out that if the sequence of weights can be written

$$\beta(t, s) = \lambda(t)\beta(t-1, s) \quad 1 \leq s \leq t-1$$

$$\beta(t, t) = 1$$

Then the estimates can be updated recursively



The Adaptive Recursive LS algorithm

$$\mathbf{R}_t = \mathbf{x}_t \mathbf{x}_t^T + \lambda(t) \mathbf{R}_{t-1}$$

$$\mathbf{h}_t = \mathbf{x}_t Y_t + \lambda(t) \mathbf{h}_{t-1}$$

$$\hat{\boldsymbol{\theta}}_t = \mathbf{R}_t^{-1} \mathbf{h}_t$$



Other formulations

1. Eliminating h_t :

$$\begin{aligned} \mathbf{R}_t &= \mathbf{x}_t \mathbf{x}_t^T + \lambda(t) \mathbf{R}_{t-1} \\ \hat{\boldsymbol{\theta}}_t &= \hat{\boldsymbol{\theta}}_{t-1} + \mathbf{R}_t^{-1} \mathbf{x}_t (Y_t - \mathbf{x}_t^T \hat{\boldsymbol{\theta}}_{t-1}) \end{aligned}$$

2. Eliminating h_t and avoiding matrix-inversion:

$$\begin{aligned} \mathbf{K}_t &= \frac{\mathbf{P}_{t-1} \mathbf{x}_t}{\lambda(t) + \mathbf{x}_t^T \mathbf{P}_{t-1} \mathbf{x}_t} \\ \hat{\boldsymbol{\theta}}_t &= \hat{\boldsymbol{\theta}}_{t-1} + \mathbf{K}_t (Y_t - \mathbf{x}_t^T \hat{\boldsymbol{\theta}}_{t-1}) \\ \mathbf{P}_t &= \frac{1}{\lambda(t)} \left(\mathbf{P}_{t-1} - \frac{\mathbf{P}_{t-1} \mathbf{x}_t \mathbf{x}_t^T \mathbf{P}_{t-1}}{\lambda(t) + \mathbf{x}_t^T \mathbf{P}_{t-1} \mathbf{x}_t} \right) \end{aligned}$$



Constant forgetting

- If $\lambda(t) = \lambda$ we call λ the *forgetting factor* and define the *memory* as

$$T_0 = \sum_{i=0}^{\infty} \lambda^i = 1 + \lambda + \lambda^2 + \lambda^3 + \lambda^4 + \dots = \frac{1}{1 - \lambda}$$

- Given a data set an optimal value of λ can be found by “trial and error”
- It is often a good idea to select the values of λ to be investigated so that the corresponding values of T_0 are approximately equidistant
- The criteria to evaluate may depend on the application, but the sum of squared one-step prediction errors is often appropriate
- The initialization period should be excluded from the evaluation



Variable forgetting

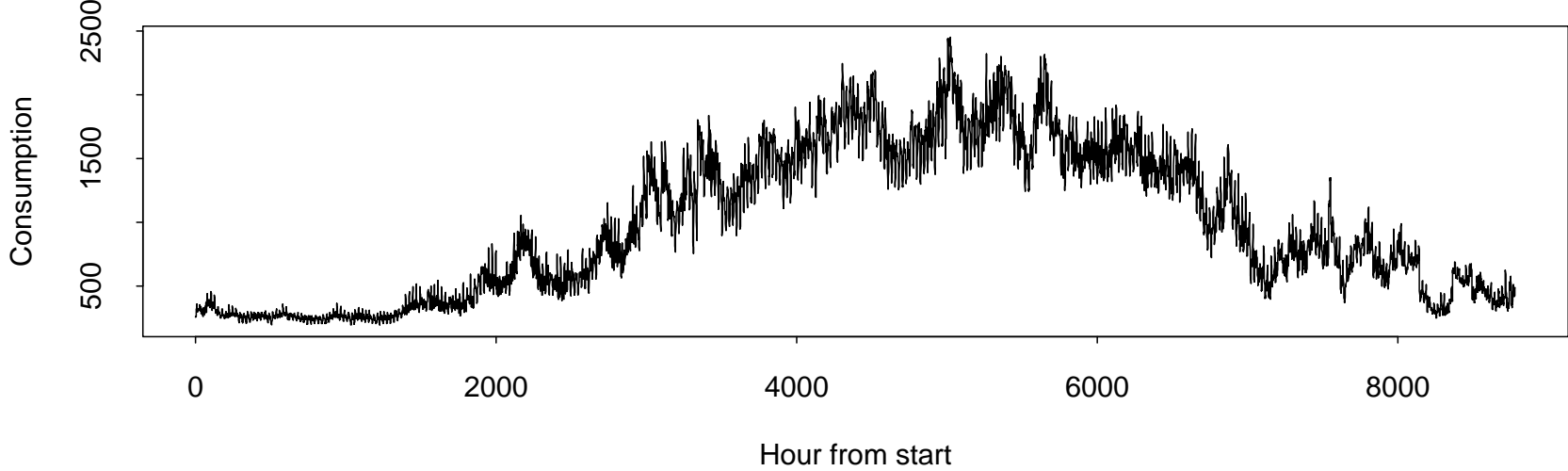
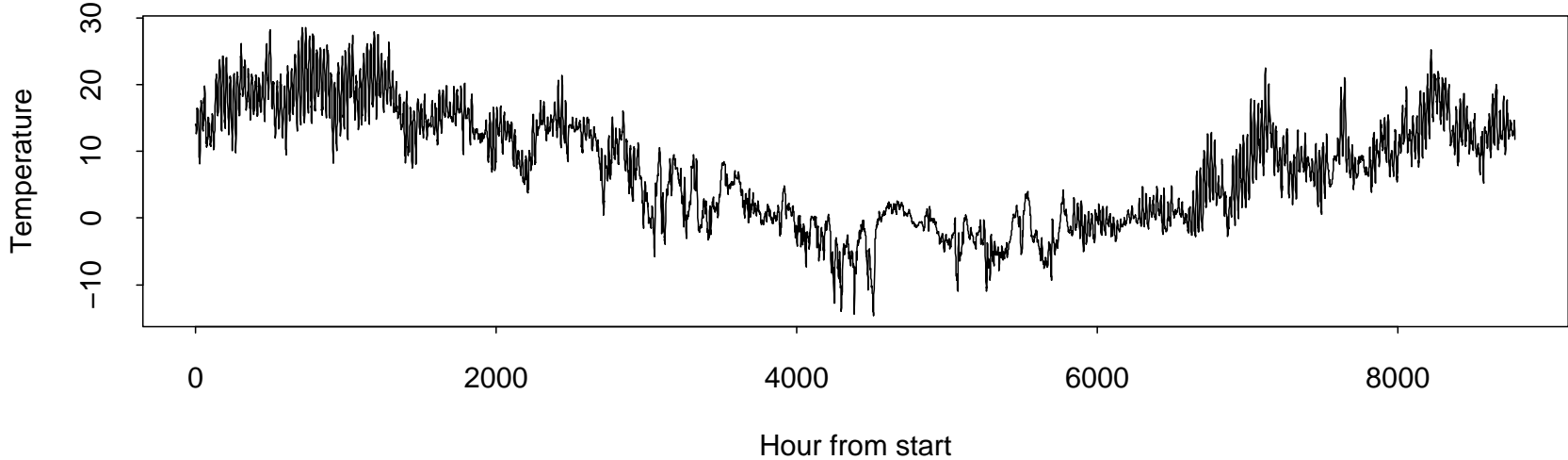
- Many methods exists
- One is based on the aim of keeping the criteria functions defining the estimates constant at S_0
- Leads to:

$$\lambda(t) = 1 - \frac{\varepsilon_t^2}{S_0 [1 + \mathbf{x}_t^T \mathbf{P}_{t-1} \mathbf{x}_t]}$$

- A lower bound λ_{min} on $\lambda(t)$ should be applied
- For optimal tuning of this method S_0 could be varied

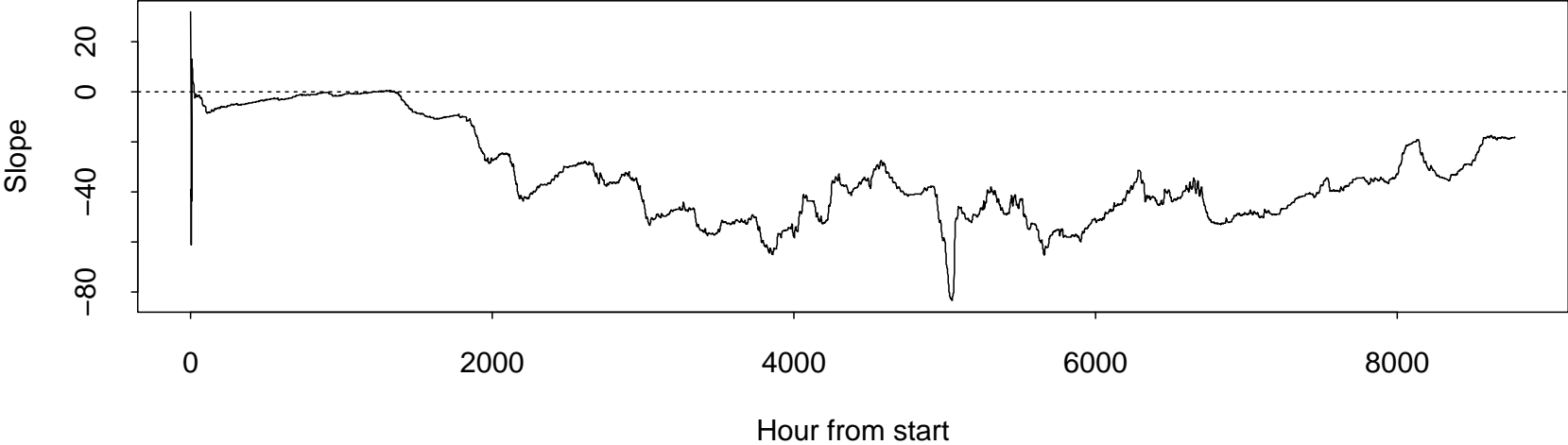
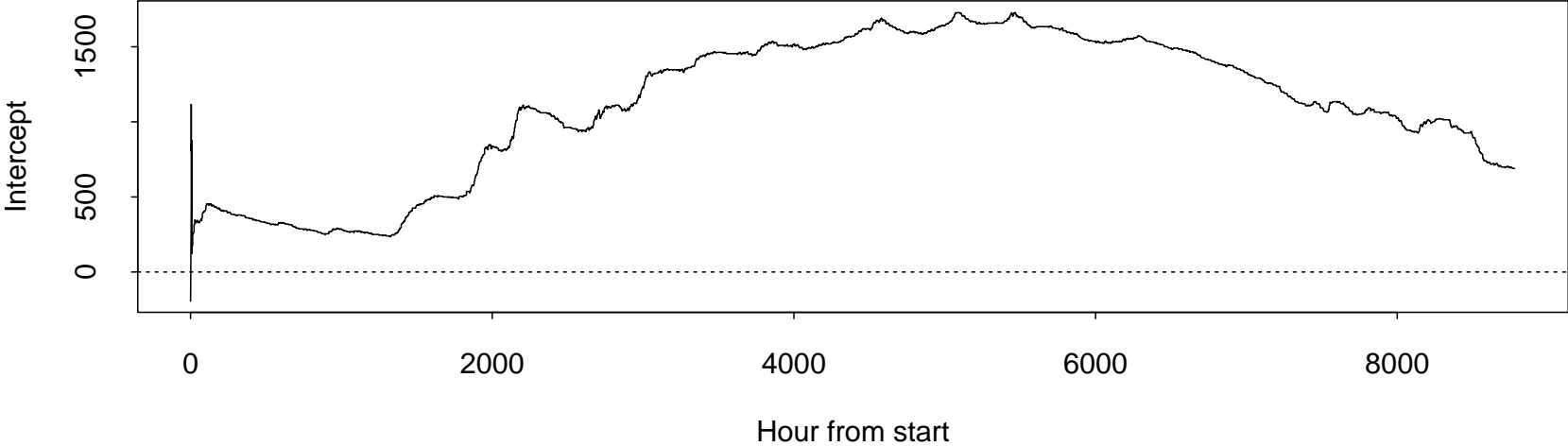


Example: $HC_t = \mu + \theta_1 T_t + \varepsilon_t$





Example: $HC_t = \mu + \theta_1 T_t + \varepsilon_t, \lambda = 0.995$





Example (cont'nd)

