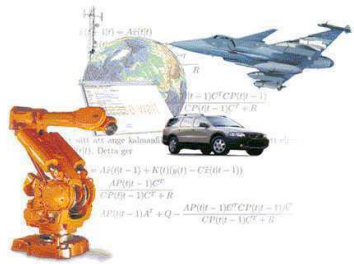


# Full Bayesian Identification of Linear Dynamic Systems Using Stable Kernels



Lennart Ljung

Reglerteknik, ISY, Linköpings Universitet  
 Joint work with  
**Gianluigi Pillonetto**, Padova University

- A new technique to identifying linear black-box dynamic systems
- Outperforms existing methods
- Further development of kernel methods using “Empirical Bayes”
- We will work essentially with Linear Regressions.

## Linear Regression for a General Linear Model

A general linear, single output model

$$y(t) = G(q, \theta)u(t) + H(q, \theta)e(t)$$

with  $nu$  inputs has an output predictor that is linear in past outputs and inputs:

$$\hat{y}(t|t-1) = \sum_{j=1}^{nu+1} \sum_{k=1}^{m_j} \tilde{h}_j(k) \tilde{u}_j(t-k)$$

where and  $\tilde{h}$  are formed from  $G$  and  $H$  and  $\tilde{u}_k = u_k, k = 1, \dots, nu$  &  $\tilde{u}_{nu+1}(k) = y(k)$ .

## A Simple Special Case: FIR model

Simple special case: A SISO FIR (finite impulse response) model :

$$y(t) = B(q)u(t) + e(t)$$

$$y(t) = b_1u(t-1) + \dots + b_nu(t-n) + e(t) = \phi(t)\theta + e(t)$$

$$Y = \Phi^T\theta + E$$

$Y$  and  $\Phi^T$  collect the data for  $t = 1 \dots N$  in column vectors and  $\theta$  is the impulse response of the system. (The general case with several inputs and an Autoregressive terms makes  $\theta$  a “segmented long” matrix, with each segment being the impulse response from each of the inputs and from the past outputs.)

## Linear Regression

Model

$$Y = \Phi^T \theta + E$$

Idea #1, LS:

$$\min_{\theta} \|Y - \Phi^T \theta\|^2 \rightarrow \hat{\theta} = [\Phi \Phi^T]^{-1} \Phi^T Y$$

Idea #2, LS with regularization:

$$\min_{\theta} \|Y - \Phi^T \theta\|^2 + \theta^T D \theta \rightarrow \hat{\theta} = [\Phi \Phi^T + D]^{-1} \Phi^T Y$$

## How to choose the regularization matrix $D$ ?

**Bayesian setting:** If a prior distribution  $\theta \in N(0, P_0)$  is given and  $E \in N(0, \sigma^2 I_N)$  then simple calculations show that the natural choice of  $D$  is

$$D = P_0^{-1} / \sigma^2$$

That makes

$$\hat{\theta} = [\Phi \Phi^T + D]^{-1} \Phi^T Y = \left( \Phi^T \Phi / \sigma^2 + P_0^{-1} \right)^{-1} \frac{\Phi^T}{\sigma^2} Y$$

which is the posterior mean  $E\theta|Y$ .

## But $P_0$ and $\sigma^2$ are not known: Hyperparameters

The hyper parameters  $\eta = \{\sigma^2, \alpha, \lambda\}$

- the noise variance  $\sigma^2$
- the parameter covariance matrix  $P_0(\alpha, \lambda)$
- $P_0(k, j) = \alpha \cdot \lambda^{\max(k, j)}$  (Size and Decay)
- Recall  $\theta$  is the impulse response from the inputs (and past outputs) to the prediction.
- So the diagonal of  $P_0$ , i.e. the elements  $\alpha \lambda^k$  give the decay of the impulse response. The correlation between two impulse response coefficients  $\tilde{h}_j$  and  $\tilde{h}_k$  decays like  $\lambda^{|j-k|}$ . It is convenient to measure both decay and correlation by the same entity  $\lambda$ , but it is easy to use different numbers for these (The “TC” and “DC” kernels, resp.)
- in the multi-input case there is one covariance matrix  $P_0$  for each input (segment of the  $\theta$  matrix)

## The ML Estimate of the Hyperparameters

Actually, the ML estimate of the hyperparameter  $\eta$  is easily found.

$$Y = \Phi^T \theta + E$$

is a sum of two Gaussian variables, so  $Y$  is Gaussian with mean zero and covariance  $P_Y(\eta) = \Phi^T P_0(\eta) \Phi + \sigma^2(\eta) I_N$ . So, given  $\eta$  the pdf of  $Y$  is known, and the ML estimate of  $\eta$  will be

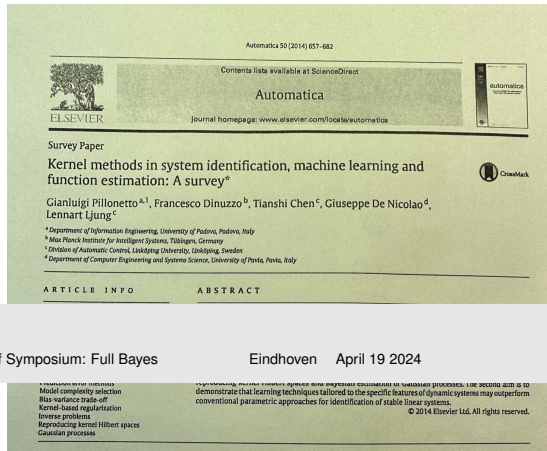
$$\hat{\eta}_{ML} = \arg \min_{\eta} Y^T P_Y(\eta)^{-1} Y + \log \det P_Y(\eta)$$

## Empirical Bayes

Inserting the ML estimate  $\hat{\eta}_{ML}$  of the hyperparameters in the regularization matrix

$$D = P_0^{-1} / \sigma^2$$

give the *Empirical Bayes* estimation method. This has proven to an effective method for identification of linear systems. Survey Paper



Lennart Ljung

Paul van der Hof Symposium: Full Bayes

Eindhoven April 19 2024

AUTOMATIC CONTROL  
REGLERTEKNIK  
LÄNKÖPINGS UNIVERSITET



## Constructing the posterior distribution by MCMC

To improve on the EB estimate we should construct the posterior distribution  $p(\eta|Y)$  more seriously by treating also the hyperparameters as random variables and

- Defining priors for the hyperparameters  $\eta$ .
- Define a chain of randomly perturbed values of  $\eta$ .
- For given values of the hyperparameters the estimate of  $\theta$  is well defined by regularized LS.
- Monitor the statistics of the corresponding sequence of  $\theta$ -values (MCMC: Markov Chain Monte Carlo - process)

Lennart Ljung

Paul van der Hof Symposium: Full Bayes

Eindhoven April 19 2024

AUTOMATIC CONTROL  
REGLERTEKNIK  
LÄNKÖPINGS UNIVERSITET



## Beyond Empirical Bayes

We are looking for  $\hat{\theta} = E(\theta|Y) = \int \theta \cdot p(\theta|Y)$  the mean of the posterior distribution. The posterior distribution of  $\theta$ :

$$p(\theta|Y) = \int p(\theta, \eta|Y) d\eta = \int p(\theta|\eta, Y) p(\eta|Y) d\eta$$

If the posterior distribution  $p(\eta|Y)$  is concentrated around the ML estimate  $\hat{\eta}_{ML}$  we can use

$$p(\theta|Y) = \int p(\theta, \eta|Y) d\eta = \int p(\theta|\eta, Y) p(\eta|Y) d\eta \approx p(\theta|\hat{\eta}_{ML}, Y)$$

which is **Empirical Bayes (EB)**.

Lennart Ljung

Paul van der Hof Symposium: Full Bayes

Eindhoven April 19 2024

AUTOMATIC CONTROL  
REGLERTEKNIK  
LÄNKÖPINGS UNIVERSITET



## Hyperparameter Priors for $\sigma^2, \alpha, \gamma$

- $\sigma^2, \alpha$ , are positive real numbers with *Jeffrey's prior*
  - "noninformative prior", its posteriors sample with the *inverse gamma distribution*  $I_g(a, b)$  (with  $a, b$  given from observations of data)
  - $I_g(a, b)$  pdf  $p(x) \sim x^{-a-1} e^{b/x}$
- $\lambda$  positive scalar less than 1.

Lennart Ljung

Paul van der Hof Symposium: Full Bayes

Eindhoven April 19 2024

AUTOMATIC CONTROL  
REGLERTEKNIK  
LÄNKÖPINGS UNIVERSITET



## The MCMC Steps: Hyperparameters $\sigma^2, \lambda, \alpha$

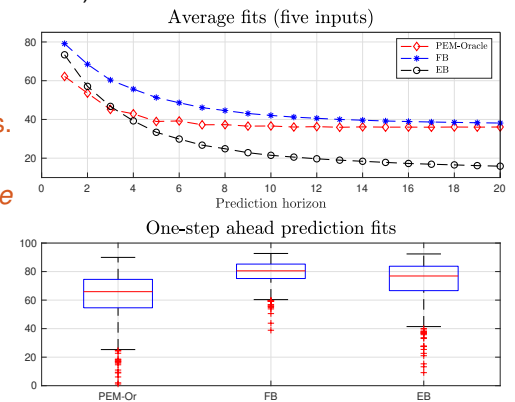
1. Initialize:  $\theta^0$  is the LS estimate,  $\sigma_0^2 =$  the loss function for the LS estimate. Set e.g.  $\lambda = 0.95, \alpha = 1$ . Set  $\Delta = 0.02$
2. For  $i = 1, 2, \dots, M$  ( $M \sim 10000$ ) repeat
  - resample the value for  $\sigma_i^2$  as the loss for the current model  
 $\sigma_i^2 = \|Y - \Phi\theta^{i-1}\|^2 / (N - 1)$   
 $[\sigma_i^2 \sim I_g(N/2, \|Y - \Phi\theta^{i-1}\|^2/2)]$
  - similarly, resample  $\alpha$  as a value drawn from an inverse gamma distribution with parameters given by the current estimate of  $\theta^{i-1}$
  - make a random change of size  $\Delta$  in  $\lambda$ . Accept the change with a probability that depends on the data fit with new and old  $\lambda$  values
  - Compute the posterior distribution of  $\theta$  for the new hyperparameters and  $\mu^i$  be mean of this distribution, while the new estimate  $\theta^i$  be drawn from this distribution.
3. The resulting estimate will be  $\hat{\theta} = \frac{1}{M} \sum_{i=1}^M \mu^i$
4. The sequence  $\{\theta^i\}$  is the sampled posterior density  $p(\theta|Y)$



## Example: Simulated Models

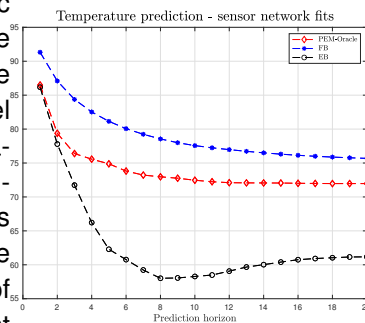
Generate randomly a system with 5 inputs. Estimate models with three different methods and compute and compare  $m$ -step predictions (1000 systems generated)

1. PEM-Oracle (red): The standard, "classical" method. Try ARMAX models of different orders. Select the order that works best (with an oracle – ideal)  $\sim 10\text{sec}$
2. EB: (black) Empirical Bayes  $< 1\text{sec}$ .
3. FB:(blue) Full Bayes (using MCMC),  $\sim 10\text{sec}$



## Real Experiment: Temperature Sensor Network

A building is covered by a network of 24 thermodynamic sensors (temperature, relative humidity, and similar). One sensor is selected to model the temperature from the measurements in the 23 other sensors. An ARMAX structure is set up with 23 inputs and one output, and the parameters of this are estimated in different ways.



## Conclusions

The "Full Bayesian Approach" has a clear edge over "Empirical Bayes" and "classical methods" and could be seen as the preferred method to estimate linear black box models. No need to deal with order/structure determination of ARX, ARMAX, Box-Jenkins models. Full story in

G. Pillonetto and L. Ljung: Full Bayesian identification of linear dynamic systems using stable kernels. Proc National Academy of Sciences, PNAS, Vol 120 (18), 2023, e2218197120



# Take it Easy — Keep on Working!



Lennart Ljung  
Paul van der Hof Symposium: Full Bayes

Eindhoven April 19 2024

AUTOMATIC CONTROL  
REGLERTEKNIK  
LINKÖPINGS UNIVERSITET

