

## Construction of an Informative Hierarchical Prior Distribution: Application to Electricity Load Forecasting

Anne Philippe  
Laboratoire de Mathématiques Jean Leray  
Université de Nantes

Workshop EDF-INRIA, 5 avril 2012, IHP - Paris

Work in collaboration with  
Tristan Launay – LMJL / EDF OSIRIS, R39  
and  
Sophie Lamarche – EDF OSIRIS, R39



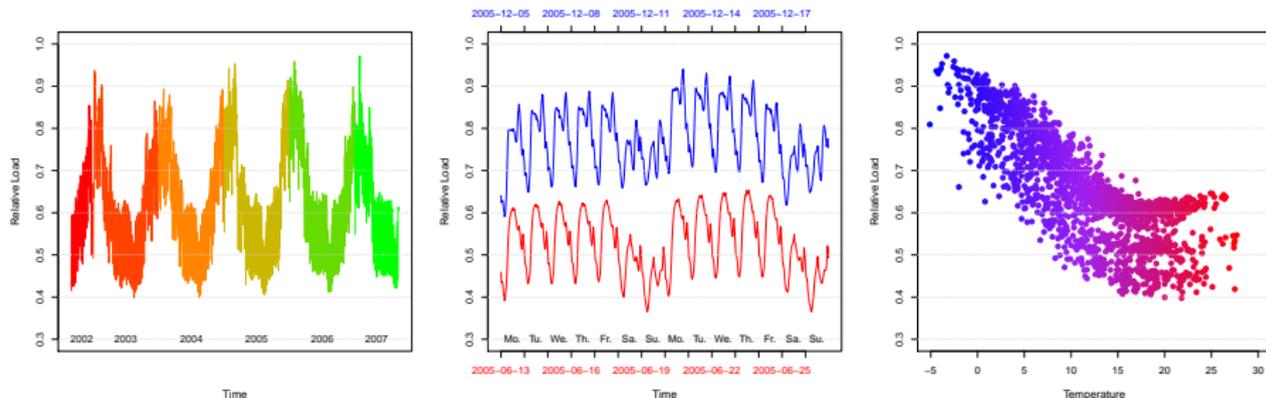
# Plan

- 1 Introduction
- 2 Bayesian approach and asymptotic results
- 3 Construction of the prior
- 4 Numerical results
- 5 Bibliography

# Plan

- 1 Introduction
- 2 Bayesian approach and asymptotic results
- 3 Construction of the prior
- 4 Numerical results
- 5 Bibliography

# The electricity load signal



## Comments

- the electricity demand exhibits multiple strong seasonalities : yearly, weekly and daily.
- daylight saving times, holidays and bank holidays induce several breaks
- non-linear link with the temperature.

# The model

The electricity load  $X_t$  is written as

$$X_t = s_t + w_t + \epsilon_t$$

- $s_t$  : seasonal part
- $w_t$  : non-linear weather-dependent part (heating)
- $(\epsilon_t)_{t \in \mathbb{N}}$  i.i.d. from  $\mathcal{N}(0, \sigma^2)$

## Heating part

$$w_t = \begin{cases} g(T_t - u) & \text{if } T_t < u \\ 0 & \text{if } T_t \geq u \end{cases}$$

- $g$  is the heating gradient,
- $u$  is the heating threshold

# Model (cont.)

## Seasonal part

$$s_t = \left[ \sum_{j=1}^p \left\{ a_j \cos \left( \frac{2\pi j}{365.25} \cdot t \right) + b_j \sin \left( \frac{2\pi j}{365.25} \cdot t \right) \right\} + \sum_{j=1}^q \omega_j \mathbb{1}_{\Omega_j}(t) \right] \cdot \sum_{j=1}^r \kappa_j \mathbb{1}_{K_j}(t)$$

- the average seasonal behaviour with a truncated Fourier series
- gaps (parameters  $\omega_j \in \mathbb{R}$ ) represent the average levels over predetermined periods given by a partition  $(\Omega_j)_{j \in \{1, \dots, d_{12}\}}$  of the calendar.
- day-to-day adjustments through shapes (parameters  $\kappa_j$ ) that depends on the so-called days' types which are given by a second partition  $(K_j)_{j \in \{1, \dots, d_2\}}$  of the calendar.

# Plan

- 1 Introduction
- 2 Bayesian approach and asymptotic results**
- 3 Construction of the prior
- 4 Numerical results
- 5 Bibliography

## Bayesian principle

Let  $X_1, \dots, X_N$  be the observations.

The likelihood of the model, conditionnally to the temperatures, is :

$$l(X_1, \dots, X_N | \theta) = \prod_{t=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} [X_t - (s_t + w_t)]^2 \right\}$$

Let  $\pi(\theta)$  be the *prior distribution* on  $\theta$

The bayesian inference is based on

- the *posterior distribution* of  $\theta$

$$\pi(\theta | X_1, \dots, X_N) \propto l(X_1, \dots, X_N | \theta) \pi(\theta)$$

- The Bayes estimate (under quadratic loss) is  $\hat{\theta} = \mathbb{E}[\theta | X_1, \dots, X_N]$

## Predictive distribution

- 

$$p(X_{N+k} | X_1, \dots, X_N) = \int_{\Theta} l(X_{N+k} | \theta) \pi(\theta | X_1, \dots, X_N) d\theta$$

- The optimal prediction (under quadratic loss) is  $\hat{X}_{N+k} = \mathbb{E}[X_{N+k} | X_1, \dots, X_N]$

# Asymptotic results for the heating part

The observations  $X_{1:n} = (X_1, \dots, X_n)$  depend on an exogenous variable  $T_{1:n} = (T_1, \dots, T_n)$  via the model

$$X_i = \mu(\eta, T_i) + \xi_i := g(T_i - u)\mathbb{1}_{[T_i, +\infty[}(u) + \xi_i,$$

for  $i = 1, \dots, n$ .

- $(\xi_i)_{i \in \mathbb{N}}$  is a sequence of i.i.d. from  $\mathcal{N}(0, \sigma^2)$ .
- $\theta_0$  will denote the true value of  $\theta$ .

$$\theta = (g, u, \sigma^2) \in \Theta = \mathbb{R}^* \times ]\underline{u}, \bar{u}[ \times \mathbb{R}_+^*.$$

# Consistency

## Assumptions (A).

- 1 There exists  $K \subset \Theta$  a compact subset of the parameter space  $\Theta$  such that the MLE  $\hat{\theta}_n \in K$  for any  $n$  large enough.
- 2 Let  $F$  be a cumulative distribution function which is continuously differentiable and  $F'$  does not vanish over  $[\underline{u}, \bar{u}]$

$$F_n(u) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[T_i, +\infty[}(u), \xrightarrow{n \rightarrow +\infty} F(u)$$

## Theorem ( Launay et al [2])

Let  $\pi(\cdot)$  be a prior distribution on  $\theta$ , continuous and positive on a neighbourhood of  $\theta_0$  and let  $U$  be a neighbourhood of  $\theta_0$ , then under Assumptions (A), as  $n \rightarrow +\infty$ ,

$$\int_U \pi(\theta | X_{1:n}) d\theta \xrightarrow{a.s.} 1.$$

# Asymptotic normality

## Theorem (Launay et al [2])

Let  $\pi(\cdot)$  be a prior distribution on  $\theta$ , continuous and positive at  $\theta_0$ , and let  $k_0 \in \mathbb{N}$  such that

$$\int_{\Theta} \|\theta\|^{k_0} \pi(\theta) \, d\theta < +\infty,$$

and denote

$$t = n^{\frac{1}{2}}(\theta - \hat{\theta}_n),$$

and  $\tilde{\pi}_n(\cdot|X_{1:n})$  the posterior density of  $t$  given  $X_{1:n}$ , then under Assumptions (A), for any  $0 \leq k \leq k_0$ , as  $n \rightarrow +\infty$ ,

$$\int_{\mathbb{R}^3} \|t\|^k \left| \tilde{\pi}_n(t|X_{1:n}) - (2\pi)^{-\frac{3}{2}} |I(\theta_0)|^{\frac{1}{2}} e^{-\frac{1}{2}t'I(\theta_0)t} \right| \, dt \xrightarrow{\mathbb{P}} 0,$$

where  $I(\theta)$  is the asymptotic Fisher Information matrix.

# Plan

- 1 Introduction
- 2 Bayesian approach and asymptotic results
- 3 Construction of the prior**
- 4 Numerical results
- 5 Bibliography

## Objective

Consider the general model  $X_t = s_t + w_t + \epsilon_t$

### Aim

the estimation and the predictions of the parametric model over a **short** dataset denoted  $B$

### Problem

- The high dimensionality of a model leads to an overfitting situation, and so the errors in prediction are larger
- the parameter of the model

$$\theta = (\eta, \sigma^2), \quad \eta = \underbrace{(a_1, \dots, a_p, b_1, \dots, b_p)}_{\text{Fourier}}, \underbrace{(\omega_1, \dots, \omega_k)}_{\text{offsets}}, \underbrace{(\kappa_1, \dots, \kappa_r)}_{\text{shapes}}, \underbrace{(g, u)}_{\text{heating part}}$$

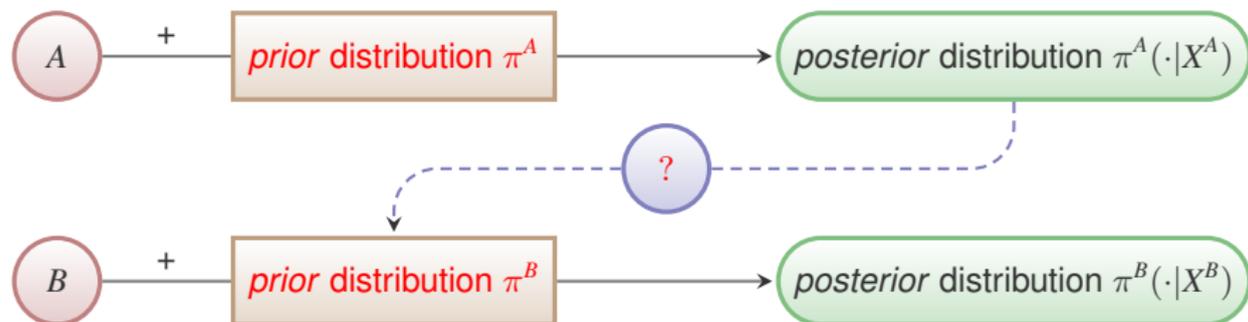
### Prior information

- A a long dataset known to share some common features with  $B$

We wish to improve parameter estimations and model predictions over  $B$  with the help of  $A$ .

## Choice of the prior distribution

Construction of the prior distribution for  $B$  the short dataset using the information coming from the long dataset  $A$



## Methodology

- we assume  $\pi^B$  belongs to a parametric family  $\mathcal{F} = \{\pi_\lambda, \lambda \in \Lambda\}$
- we want to pick  $\lambda^B$  such that  $\pi^B$  is «close» to  $\pi^A(\cdot|X^A)$

Let  $T : \mathcal{F} \rightarrow \Lambda$  be an operator such that for every  $\lambda \in \Lambda$ ,

$$T[\pi_\lambda] = \lambda$$

### Choice of the *prior* distribution on $B$

- we select  $\lambda^B$  «proportional to»  $T[\pi^A(\cdot|y^A)]$  in the sense that

$$\lambda^B = K T[\pi^A(\cdot|y^A)]$$

where  $K : \Lambda \rightarrow \Lambda$  is a diagonal operator

$$K = \text{Diag}(k_1, \dots, k_j, \dots)$$

- the operator  $K$  can be interpreted as a similarity operator between  $A$  and  $B$  (the diagonal components of  $K$  are hyperparameters of the prior)
- we give  $k_j$  a vague hierarchical prior distribution centred around  $q_j$ ,
- the prior on  $q_j$  is vague and centred around 1.

## Methodology (cont.)

### Example 1 : Method of Moments

We assume that the elements of  $\mathcal{F}$  can be identified via their  $m$  first moments.

$$T[\pi_\lambda] = F(\mathbb{E}(\theta), \dots, \mathbb{E}(\theta^m)) = \lambda$$

The expression of  $\lambda^B$  then becomes

$$\lambda^B = \mathbf{K}F(\mathbb{E}(\theta|X^A), \dots, \mathbb{E}(\theta^m|X^A))$$

### Example 2 : Conjugacy

We assume that  $\mathcal{F}$  is the family of priors conjugated for the model.

$$\pi^A \in \mathcal{F} \Rightarrow \pi^A(\cdot|X^A) \in \mathcal{F}$$

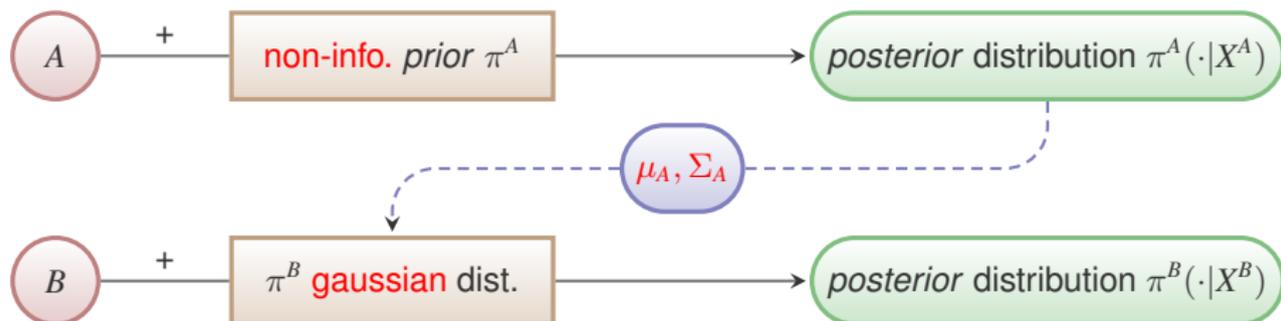
let  $\lambda^A(X^A) \in \Lambda$  be the parameter of *posterior* distribution  $\pi^A(\cdot|X^A)$

The expression of  $\lambda^B$  thus reduces to

$$\lambda^B = \mathbf{K}\lambda^A(X^A)$$

# Application to the model of electricity load

- Non informative prior for the the long dataset A
- $\pi^B \in \mathcal{F} = \{\text{gaussian distributions}\}$



non-informative *prior* for A

- we consider *prior* distribution of the form :  $\pi^A(\theta) = \pi(\eta)\pi(\sigma^2)$

- ★  $\pi(\eta) \propto 1$

- ★  $\pi(\sigma^2) \propto \sigma^{-2}$

Jeffreys prior for a standard regression model

*posterior* distribution

- the *posterior* distribution associated is well defined

$$\int_{\Theta} l(X_1, \dots, X_N | \theta) \pi^A(\theta) d\theta < +\infty$$

- Explicit forms of  $\mu^A, \Sigma^A$  are not available  $\rightsquigarrow$  MCMC algorithm

The hierarchical prior that we use is built as follows :

- the *prior* distribution is of the form :

$$\pi(\theta, k, l, q, r) = \pi(\eta|k, l)\pi(k|q, r)\pi(l)\pi(q)\pi(r)\pi(\sigma^2)$$

- prior* : the parameter of  $B$  are close to  $\mu^A$  and  $\Sigma^A$

$$\star \pi(\sigma^2) \propto \sigma^{-2}$$

$$\star \eta|k, l \sim \mathcal{N}(M^A k, l^{-1} \Sigma^A)$$

$$M^A = \text{Diag } \mu_A$$

- hyperparameters  $k$  and  $l$  correspond to the similarity between  $A$  and  $B$

$$\star k_j|q, r \sim \mathcal{N}(q, r^{-1}), l \sim \mathcal{G}(a_l, b_l)$$

$$a_l = b_l = 10^{-3}$$

- hyperparameters  $q$  and  $r$  are more general indicators of how close  $A$  and  $B$  are,  
 $q$  corresponding to the mean of the coordinates of  $k$   
 $r$  is the inverse-variance of the coordinates of  $k$

$$\star q \sim \mathcal{N}(1, \sigma_q^2), r \sim \mathcal{G}(a_r, b_r)$$

$$\sigma_q^2 = 10^4, a_r = b_r = 10^{-6}$$

# Plan

- 1 Introduction
- 2 Bayesian approach and asymptotic results
- 3 Construction of the prior
- 4 Numerical results**
- 5 Bibliography

# Simulation

## Dataset A.

- We simulated 4 years of daily data for  $A$  with parameters chosen to mimic the typical electricity load of France .
  - ★ 4 frequencies used for the truncated Fourier series
  - ★ 7 daytypes : one daytype for each day of the week
  - ★ 2 offsets to simulate the daylight saving time effect
- The temperatures ( $T_i$ ) are those measured from September 1996 to August 2000 at 10 :00AM.

Dataset B. We simulated 1 year of daily data for  $B$  with parameters :

$$\begin{aligned}
 \text{seasonal : } \alpha_i^B &= k\alpha_i^A, & \alpha &= (a, b, \omega) & \forall i &= 1, \dots, d_\alpha \\
 \text{shape : } \kappa_j^B &= \kappa_j^A, & & & \forall j &= 1, \dots, d_\beta \\
 \text{heating : } g^B &= kg^A, & u^B &= u^A. & & 
 \end{aligned}$$

We compare the quality of the bayesian predictions associated to

- the hierarchical prior
- the non-informative prior.

We simulate an extra year of daily data  $B$  to evaluate the prediction

# The comparison criterion

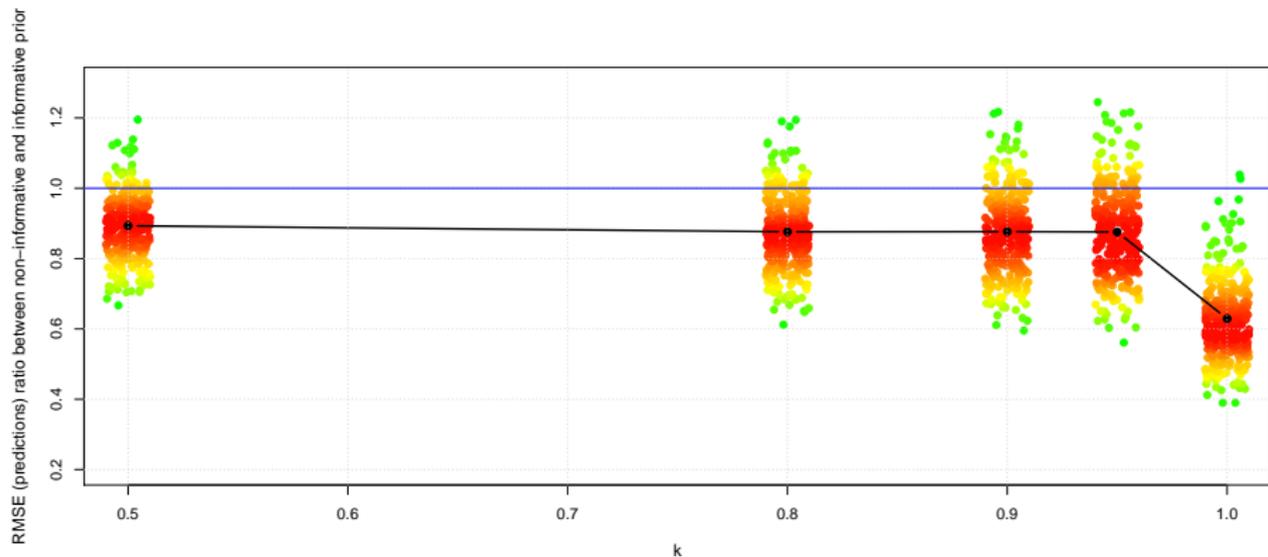
- let  $X_t = f_t(\theta) + \epsilon_t$ , the prediction error can be written as

$$X_t - \widehat{X}_t = \underbrace{[X_t - f_t(\theta)]}_{\text{noise } \epsilon_t} + \underbrace{[f_t(\theta) - \widehat{X}_t]}_{\text{difference w.r.t. the true model}}$$

- given that we want to validate our model on simulated data, we choose to consider the quadratic distance between the real and the predicted model over a year as our quality criterion for a model, i.e. :

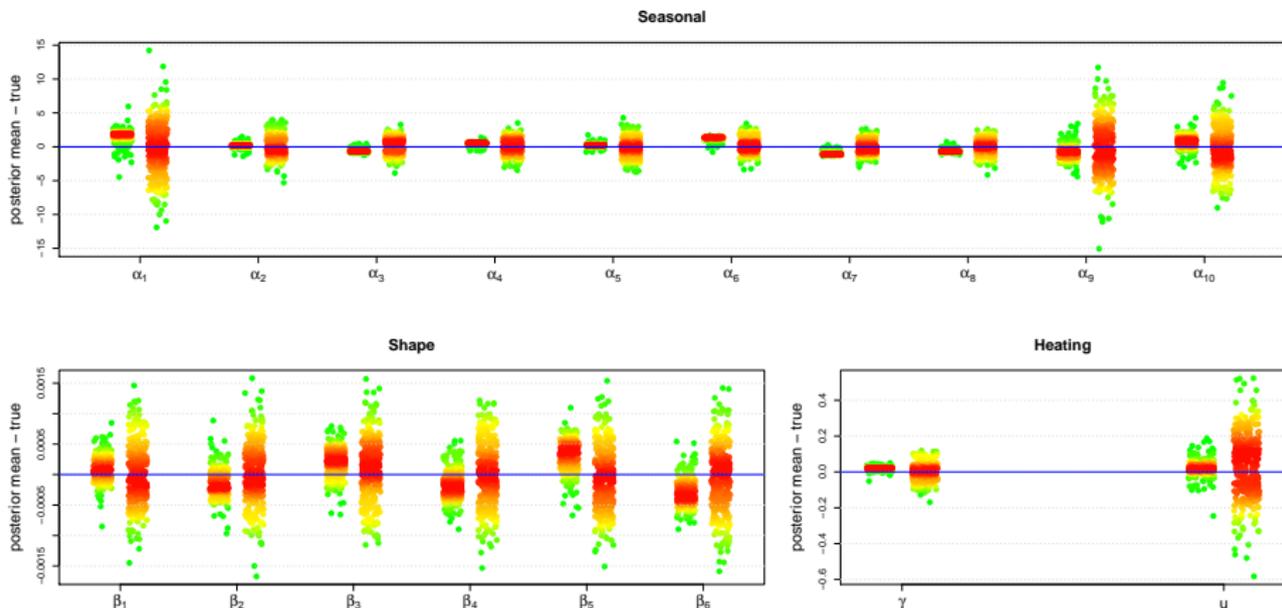
$$\sqrt{\frac{1}{365} \sum_{t=1}^{365} [f_t(\theta) - \widehat{X}_t]^2}.$$

## Comparison : RMSE info. / RMSE non-info.



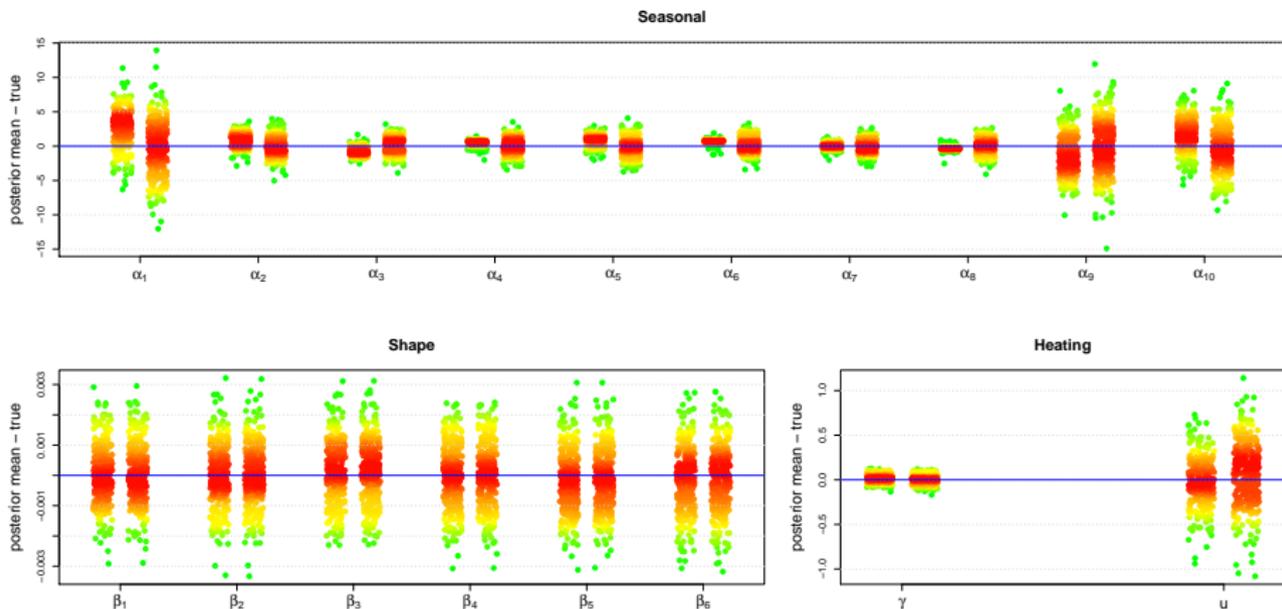
- abscissa :  $k$ , similarity coefficient between seasonality and heating gradient of populations  $A$  and  $B$  (other parameters remain untouched)
- ordinate : RMSE info. / RMSE non-info.
  - ★ whether similarity between  $A$  and  $B$  is ideal ( $k = 1$ ) or not ( $k \neq 1$ ), *prior information improves predictions' quality*

## Parameters' estimation : info. vs non-info.

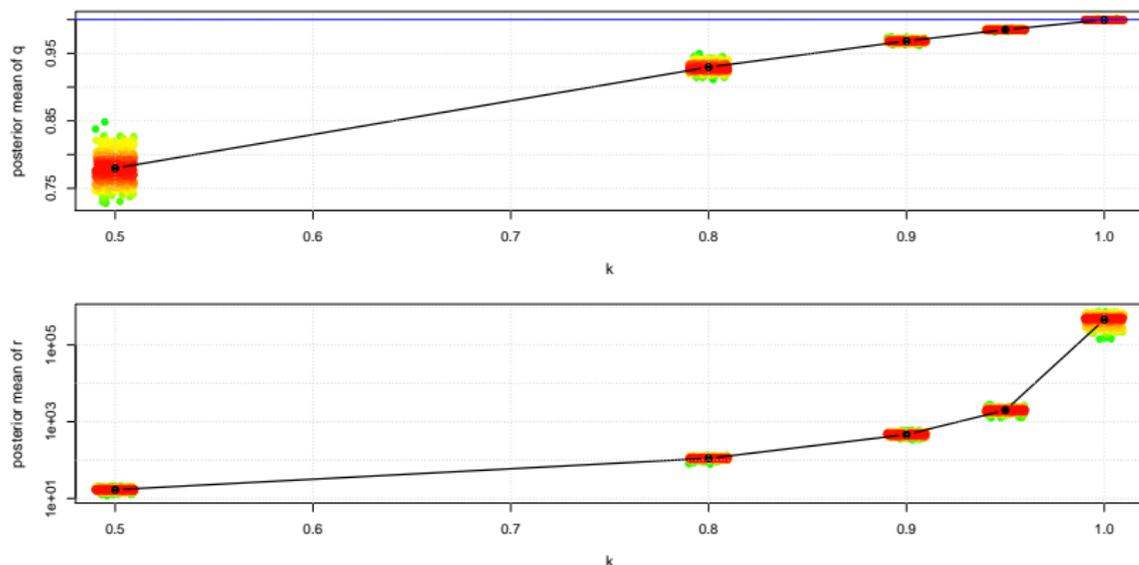


- 300 replications for the ideal situation ( $k = 1$ )
- ordinate :  $\hat{\theta} - \theta$  for the informative *prior* (left) and the non-informative *prior* (right)
  - ★ *prior* information improves the estimations a lot

## Parameters' estimation : info. vs non-info. (cont.)



- 300 replications for a difficult situation ( $k = 0.5$ )
- ordinate :  $\hat{\theta} - \theta$  for the informative *prior* (left) and the non-informative *prior* (right)
  - ★ *prior* information improves the estimations only slightly

Hyperparameters' estimation ( $q$  and  $r$ ) :  $k_j \sim \mathcal{N}(q, r^{-1})$ 

- abscissa :  $k$ , similarity coefficient between seasonality and heating gradient of populations  $A$  and  $B$  (other parameters remain untouched)
- ordinate : Bayes estimator for  $q$  and  $r$  for the  $5 \times 300$  replications considered
  - ★  $q$  et  $r$  correspond to the mean and precision (inverse of variance) of the similarity coefficients  $k_j$

# Application

## Dataset

- For  $A$  the dataset corresponds to a specific population in France frequently referred to as “non-metered”
- $B_1$  is a subpopulation of  $A$
- $B_2$  covers the same people that  $A$  (ERDF)
- The sizes (in days) of the datasets are
  - ★ population  $A$  : non metered (833 days)
  - ★ population  $B_1$  : res1+2 (177 +30 days)
  - ★ population  $B_2$  : non metered ERDF (121 +30 days)

## Dataset

- ★ 4 frequencies used for the truncated Fourier series
- ★ 5 daytypes
- ★ 2 offsets

We kept the last 30 days of each  $B$  out of the estimation datasets and assessed the model quality over the predictions for those 30 days.

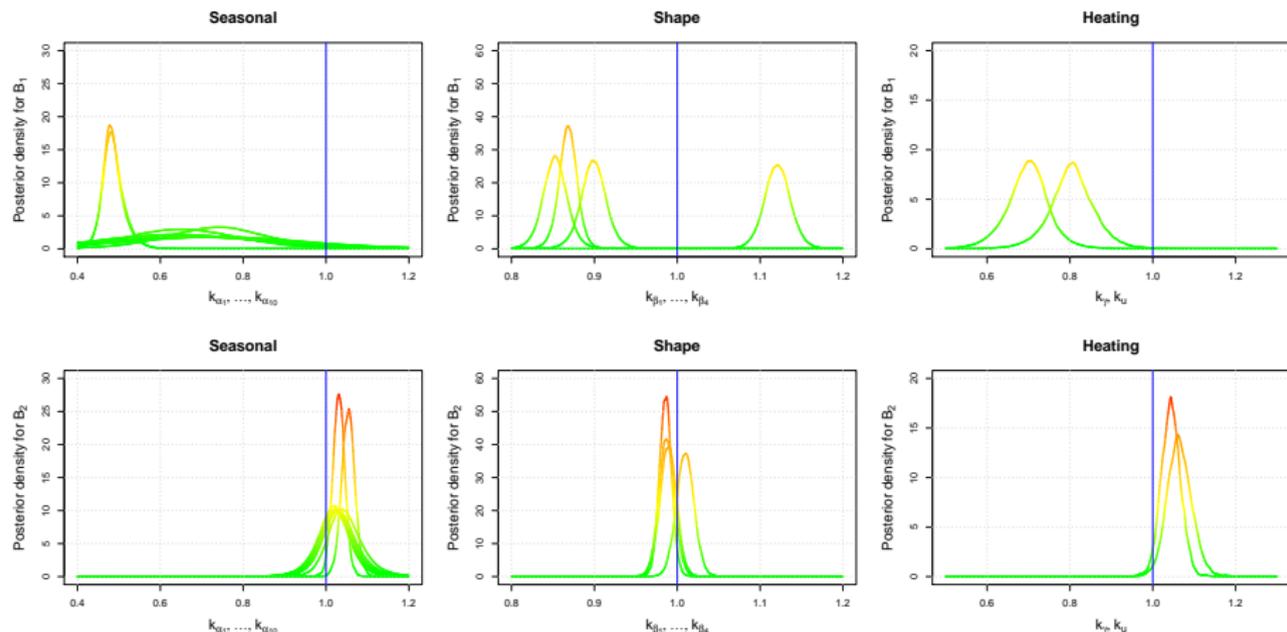
## Results

$B_1$	non-informative	hierarchical	comparison
RMSE est.	775.93	786.97	+1.42%
RMSE pred.	1863.25	894.00	-52.01%
MAPE est.	4.00	3.93	-0.07
MAPE pred.	19.37	9.30	-10.07

$B_2$	non-informative	hierarchical	comparison
RMSE est.	1127.60	1202.32	+6.62%
RMSE pred.	2286.42	1339.14	-41.83%
MAPE est.	2.82	2.98	+0.15
MAPE pred.	8.65	3.48	-5.17

**TABLE:** Results for the dataset  $B_1$  (top) and  $B_2$  (bottom). RMSE is the “root mean square error” and MAPE is the “mean absolute percentage error”. Both of these common measures of accuracy were computed on the estimation (est.) and prediction (pred.) parts of the two datasets.

Estimated posterior marginal distributions of  $k_i$  for the hierarchical method and for both datasets  $B_1$  (upper row) and  $B_2$  (lower row).



Estimation of hyperparameters  $(\hat{q}, \hat{r}) = \begin{cases} (0.73, 24.48) & \text{on } B_1 \\ (1.02, 795.16) & \text{on } B_2 \end{cases}$

# Plan

- 1 Introduction
- 2 Bayesian approach and asymptotic results
- 3 Construction of the prior
- 4 Numerical results
- 5 Bibliography**

## References I



T. Launay, A. Philippe and S. Lamarche (2011) Construction of an informative hierarchical prior distribution. Application to electricity load forecasting.

<http://arxiv.org/abs/1109.4533v2>



T. Launay, A. Philippe and S. Lamarche (2012) Consistency of the posterior distribution and MLE for piecewise linear regression

<http://arxiv.org/abs/1203.4753>