

# Quand les données sont des courbes

*Philippe BESSE*

Laboratoire de Statistique et Probabilités

UMR CNRS 5583

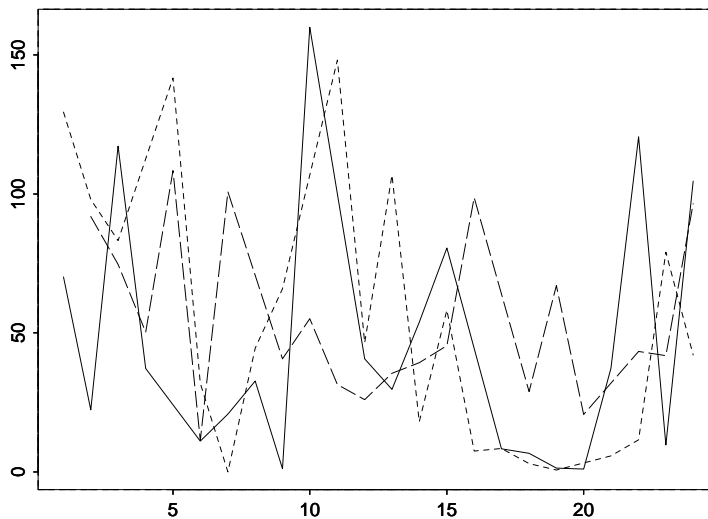
Université Paul Sabatier Toulouse III

`besse@math.ups-tlse.fr`

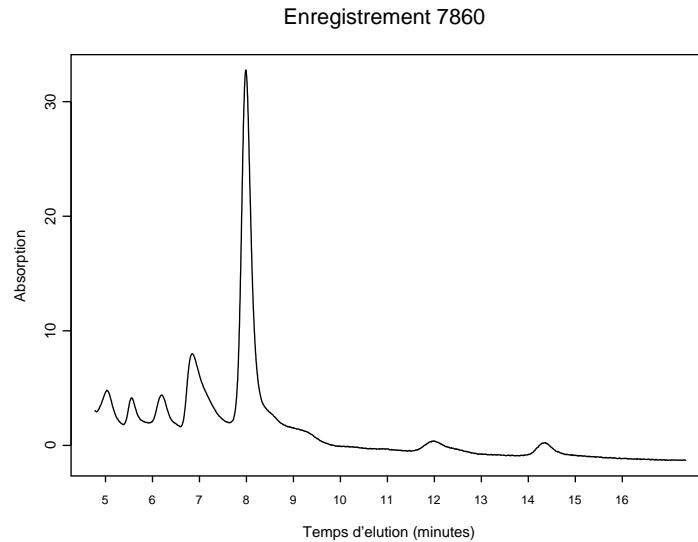
`www.lsp.ups-tlse.fr/Besse`

# 1 Introduction

- Problèmes naturellement de plus en plus fréquents.
- Historique : météorologie, chimométrie, traitement du signal...
- Bibliographie : Deville (1974), Dauxois et Pousé (1976), Besse et Ramsay (1986), Ramsay et Silverman (1997) ...
- Spécificité 1. nature fonctionnelle ou non des données,
  - base de représentation (Fourier, splines, ondelettes, fonctions propres),
  - lissage des courbes ou débruitage du signal,
  - complexité des outils mathématiques
- Spécificité 2. très forte dimensionnalité.
- Exploration (ACP) puis modélisation.
- Intérêt de la démarche sur des exemples.



*Trois exemples de courbes décrivant la pluviométrie mensuelle durant 2 ans.*



*Chromatogramme d'un jus d'orange.*

## 2 Exploration

### 2.1 ACP de courbes bruitées

#### 2.1.1 Données

- $n$  trajectoires  $z_i$ , observées en  $p$  d'instant  $t_1, \dots, t_p$
- $n$  répétitions indépendantes d'un modèle de régression non-paramétrique

$$x_j = z(t_j) + \varepsilon_j ; E(\varepsilon_j) = 0, E(\varepsilon_j \varepsilon_k) = \sigma^2 \delta_{jk}, j, k = 1, \dots, p$$
$$a \leq t_1 < t_2 < \dots < t_p \leq b$$

## 2.1.2 Modèle et estimation

- Estimation **simultanée** de  $n$  régressions non paramétriques
- **contraintes** de régularité **et** de dimension :

$$\mathbf{x}_i = \mathbf{z}_i + \boldsymbol{\varepsilon}_i; \quad i = 1, \dots, n$$

avec

$$\left\{ \begin{array}{l} \mathbf{E}(\boldsymbol{\varepsilon}_i) = 0 \text{ et } \mathbf{E}(\boldsymbol{\varepsilon}_i \boldsymbol{\varepsilon}_i) = \sigma^2 \mathbf{I}, \\ \sigma \text{ inconnue, } (\sigma > 0) \\ \mathbf{x}_i \text{ indépendant de } \boldsymbol{\varepsilon}_{i'}, \quad i' = 1, \dots, n, \\ \mathbf{x}_i \in A_q \text{ p.s. et } \|\mathbf{x}_i\|_m^2 \leq c \text{ p.s..} \end{array} \right.$$

$$\min_{\mathbf{z}_i, A_q} \left\{ \sum_{i=1}^n w_i \left( \|\mathbf{z}_i - \mathbf{x}_i\|_{\mathbf{I}}^2 + \ell \|\mathbf{z}_i\|_{\mathbf{M}}^2 \right) ; \mathbf{z} \in A_q, \dim A_q = q \right\}$$

**PROPOSITION 1.** — *La solution du problème est donnée par :*

$$\hat{\mathbf{z}}_i = \mathbf{A}_\ell^{1/2} \hat{\mathbf{P}}_q \mathbf{A}_\ell^{1/2} \mathbf{x}_i + \mathbf{A}_\ell \bar{\mathbf{x}}, \quad i = 1, \dots, n.$$

*La matrice  $\hat{\mathbf{P}}_q = \mathbf{V}_q \mathbf{V}'_q$  est la projection orthogonale sur le sous-espace  $\hat{\mathbf{E}}_q$  engendré par les  $q$  vecteurs propres de la matrice*

$$\mathbf{A}_\ell^{1/2} \mathbf{S} \mathbf{A}_\ell^{1/2}.$$

*associés aux  $q$  plus grandes valeurs propres.*

- Interpolation spline des valeurs de  $\hat{\mathbf{z}}_i$ .

### Solution équivalente :

- D.V.S. de  $(\mathbf{X}\mathbf{A}_\ell^{1/2}, \mathbf{I}, \mathbf{D})$
- donc diagonalisation de  $\mathbf{A}_\ell^{1/2}\mathbf{S}\mathbf{A}_\ell^{1/2}$ .
- Trajectoires projetées sur Vect.

$$\tilde{\mathbf{v}}_j = \mathbf{A}_\ell^{1/2}\mathbf{v}_j, \quad j = 1, \dots, q.$$

- Trajectoires estimées  $\hat{\mathbf{z}}_i$  se décomposent de manière équivalente sur la base  $\mathbf{A}_\ell^{-1}$ -orthonormée des  $\{\tilde{\mathbf{v}}_j\}$  par projection des données transformées  $\mathbf{A}_\ell\mathbf{x}_i$  :

$$\hat{\mathbf{z}}_i = \mathbf{A}_\ell\bar{\mathbf{x}} + \sum_{j=1}^q \langle \tilde{\mathbf{v}}_j, \mathbf{A}_\ell\mathbf{x}_i \rangle_{\mathbf{A}_\ell^{-1}} \tilde{\mathbf{v}}_j.$$

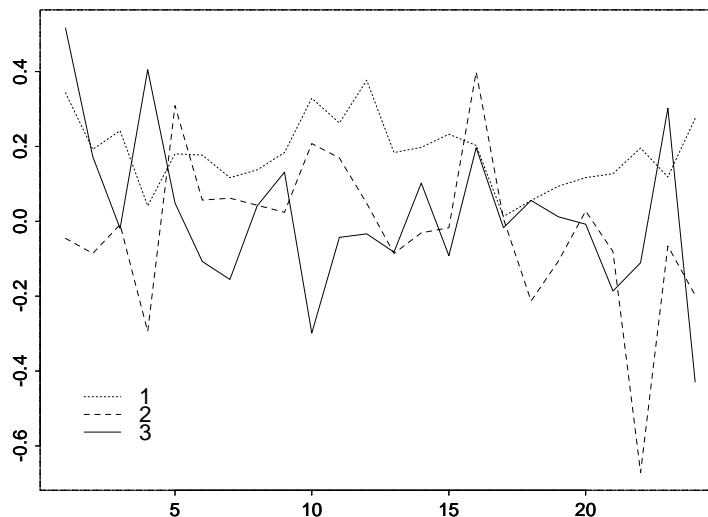
- **Dimension** et paramètre de **lissage** : optimiser un critère de **stabilité**.



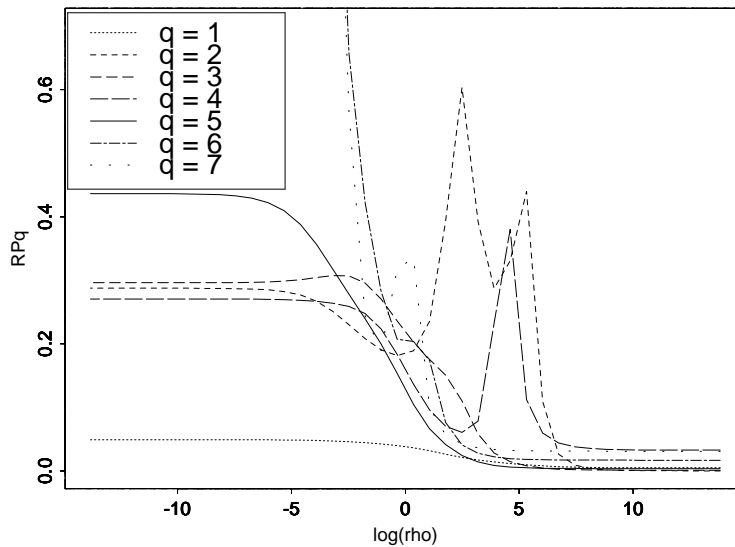
## 2.2 Exemples : ACP de séries climatiques

### 2.2.1 ACP des précipitations

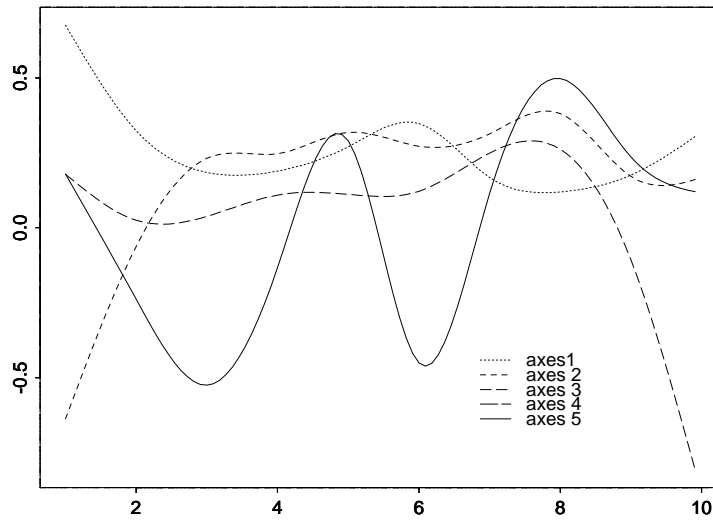
- Données : 130 courbes de précipitations mensuelles sur deux ans.



*Trois premières fonctions propres de l'ACP classique.*



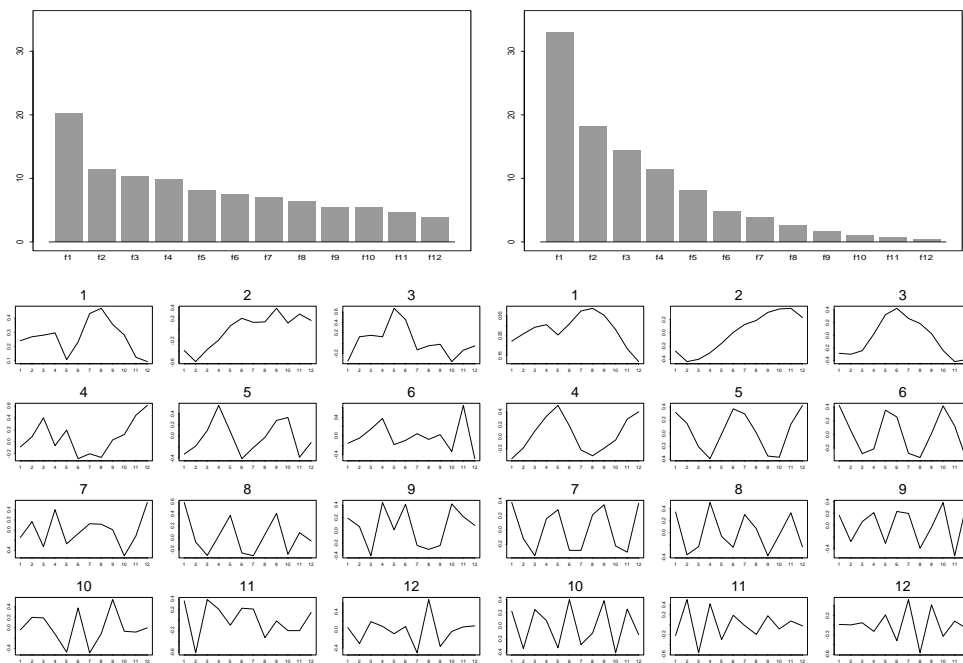
*Stabilité du sous-espace de projection.*

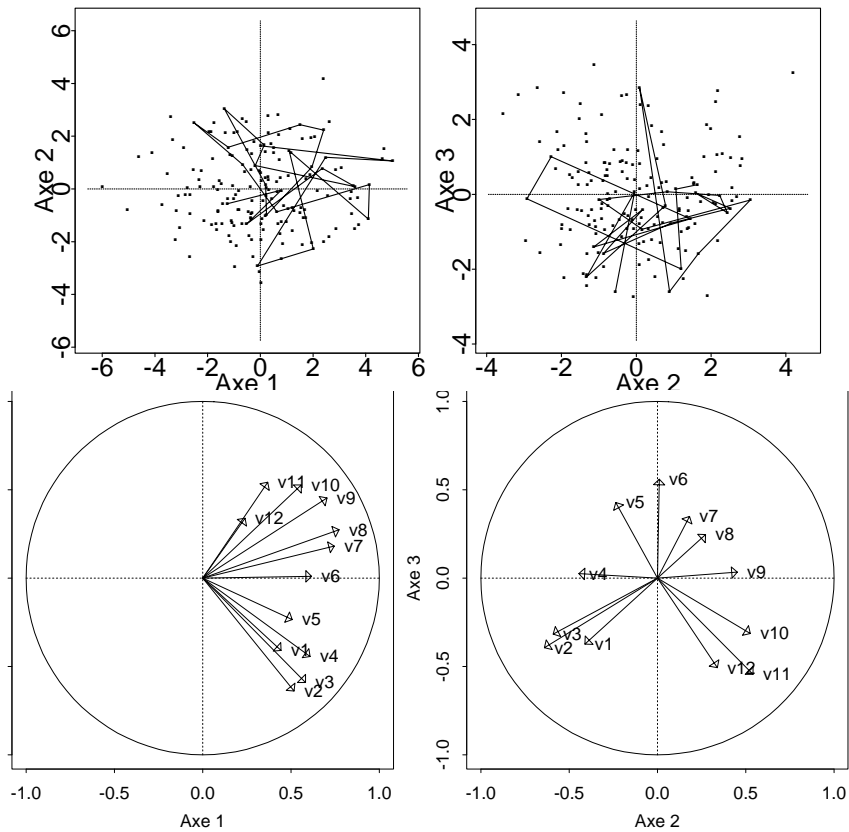


*Cinq premières composantes principales.*

## 2.2.2 ACP de températures

- Données : températures centrales en Angleterre depuis 1659.
- Moyennes mensuelles de plusieurs villes.





## 3 Modélisation de données fonctionnelles

### 3.1 Introduction

- Nature **fonctionnelle** des données comme pour **exploration**.
- **Dimensionnalité** : colinéarité, conditionnement, identifiabilité.
  - Régression sur **composantes principales** (El Niño).
  - **Régularisation** comme en *ridge*.
  - Régression **PLS**.
  - **Sélection** de fonctions de base (chromatogrammes).

## 3.2 Modèle linéaire et extensions

Ramsay et Dalzell, 1991 ; Hastie et Mallows, 1993).

soit  $(X, Y)$  un couple de variables aléatoires à valeurs dans  $H \times \mathbb{R}$

$$\mathbb{E}[Y|X = x] = \int_0^1 \psi(t)x(t) dt, \quad x \in H,$$

$\psi \in H$  coefficient de régression fonctionnel à estimer.

### 3.2.1 Modèle linéaire fonctionnel

- $\Gamma$  l'opérateur de covariance de  $X$
- $\Delta$  l'opérateur de covariance croisée de  $(X, Y)$ .

$$\Delta x = \int_0^1 \mathbb{E}[X(t)Y] x(t) dt, \quad x \in H.$$

On vérifie :

$$\Delta = \Psi\Gamma,$$

où  $\Psi(x) = \langle \psi, x \rangle$ .

**Problème** : inverser  $\Gamma$

**Solutions** :

- Dans l'espace engendré par les  $q$  premiers **vecteurs propres** (ACP).
- **Régularisation** comme en régression *ridge*.

Soit un échantillon  $(\mathbf{x}_i, y_i)$ ,  $i = 1, \dots, n$  où  $\mathbf{x}_i = (x_i(t_j), j = 1, \dots, p) \in \mathbb{R}^p$  obtenu par discrétisation de la fonction ou courbe  $x_i$ .



## 3.2.2 Régression sur composantes principales

$$\widehat{\boldsymbol{\psi}}(t) = \sum_{j=1}^q \beta_j \widehat{v}_j(t).$$

- Les composantes principales  $\langle \widehat{z}_i, \widehat{v}_j \rangle_{\mathbf{L}^2} = \mathbf{z}'_i \mathbf{N} \mathbf{v}_j$  sont les variables explicatives

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^q} \sum_{j=1}^n w_i \left( Y_i - \sum_{j=1}^q \beta_j \mathbf{z}'_i \mathbf{N} \mathbf{v}_j \right)^2.$$

### 3.2.3 Approche par pénalisation

Hastie et Mallows (1993), Marx et Eilers (1999)

$$\min_{\psi} \sum_{i=1}^n w_i (Y_i - \psi' \mathbf{N} \mathbf{x}_i)^2 + \ell \|\psi\|_2^2,$$

$$\min_{\psi} \psi' \mathbf{N} \mathbf{S} \mathbf{N} \psi - 2 \psi' \mathbf{N} \mathbf{Z}' \mathbf{D} \mathbf{Y} + \ell \psi' \mathbf{M} \psi.$$

La solution s'écrit alors

$$\hat{\psi} = (\mathbf{N} \mathbf{S} \mathbf{N} + \ell \mathbf{M})^{-1} \mathbf{N} \mathbf{Z}' \mathbf{D} \mathbf{Y}.$$

- Choix du paramètre de lissage.

### 3.2.4 Modèle linéaire général et autres extensions

Marx et Eilers (1999)

$$\mathbb{E}[Y|X = x] = g^{-1}(\langle \psi, x \rangle), \quad x \in H,$$

où  $g$  est la fonction lien de type logit, log, ...

- Vraisemblance pénalisée :

$$\max_{\psi} \sum_{i=1}^n \log \mathcal{L}(y_i, \mathbf{x}_i, \psi) - \ell \|\psi\|_2^2$$

## 3.3 Prédiction fonctionnelle

### 3.3.1 Problème

Bosq (2000), Besse et Cardot (1996) et Besse et coll. (2000).

- $(X_t)_{t \in \mathbb{R}}$  découpé en intervalles  $\delta$ ,
- $(Z_n)_{n \in \mathbb{Z}}$ , à valeurs dans  $H = L^2[0, 1]$

### 3.3.2 Modèle ARH(1)

•  $(Z_i)_{i \in \mathbb{Z}}$  processus **auto-régressif hilbertien** du premier ordre d'espérance  $a \in H$  et d'opérateur d'auto-corrélation  $\rho$  :

- $\forall i \in \mathbb{Z}, \quad Z_i - a = \rho(Z_{i-1} - a) + \epsilon_i.$
- $\sum_{n \geq 1} \|\rho^n\| < +\infty.$
- $\{\epsilon_i\}$  est centré i.i.d. dans  $H$  de variance finie,
- $\mathbb{E} \|\epsilon_i\|_H^2 = \sigma^2 < +\infty.$
- $\mathbb{E}(Z_{i+1} | Z_i, Z_{i-1}, \dots) - a = \rho(Z_i - a), \quad i \in \mathbb{Z}.$

$\Gamma = \mathbb{E}((Y_i - a) \otimes (Y_i - a))$  et  $\Delta = \mathbb{E}((Y_i - a) \otimes (Y_{i+1} - a))$  vérifient :

$$\Delta = \rho\Gamma.$$

- Inverser  $\Gamma$  pour estimer  $\rho$ .
- Régression sur **composantes principales**.
- $\hat{\rho}$  approché dans la base engendrée par les premières fonctions propres de  $\Gamma$  :

$$\hat{\rho}(s, t) = \sum_{j=1}^q \sum_{j'=1}^q \hat{\beta}_{jj'} \hat{v}_j(t) \hat{v}_{j'}(s).$$

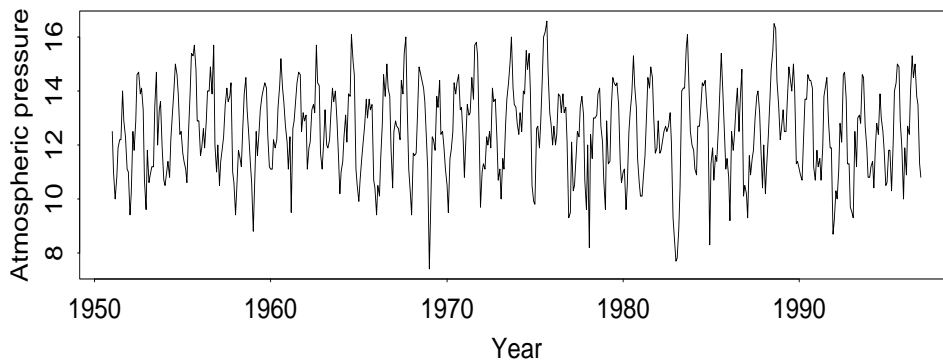
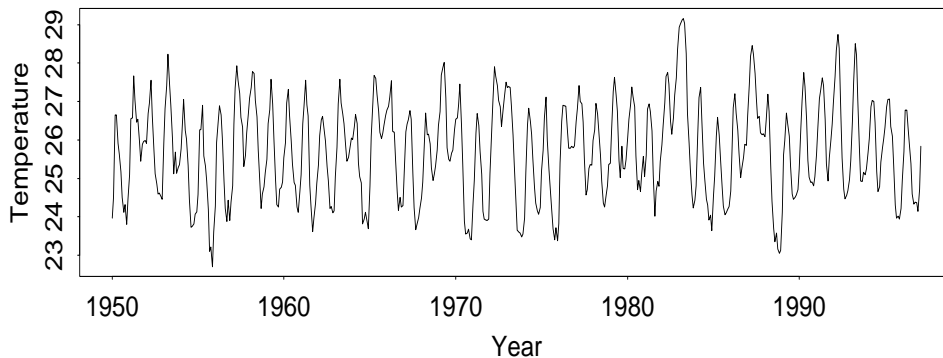
$$\min_{\hat{z}_i \in H_q} \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{p} \sum_{j=1}^p (z_i(t_j) - \hat{z}_i(t_j))^2 + \ell \|D^2 \hat{z}_i\|_{L^2}^2 \right)$$

$$\hat{\Gamma}_{q,\ell} = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i' \mathbf{N}, \quad \hat{\Delta}_{q,\ell} = \frac{1}{n-1} \sum_{i=1}^{n-1} \hat{\mathbf{x}}_{i+1} \hat{\mathbf{x}}_i' \mathbf{N}.$$

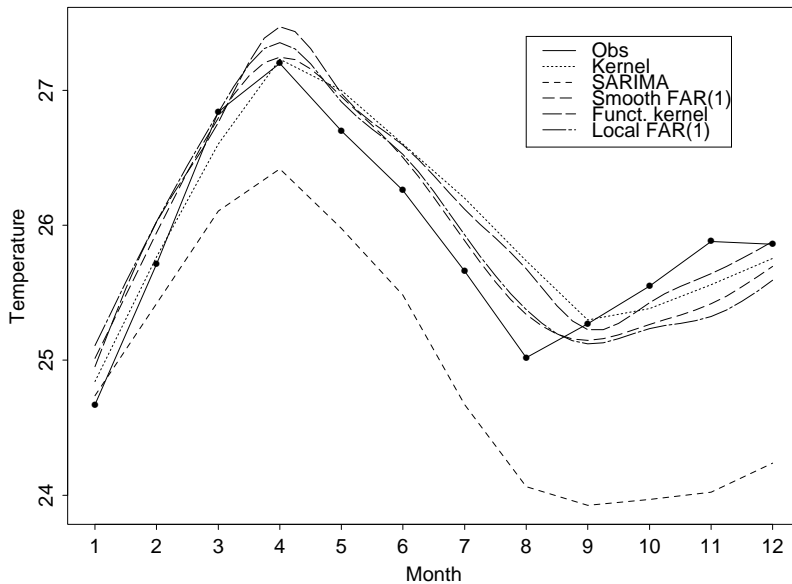
- Soit  $\hat{\Gamma}_{q,\ell}^{-1}$ , l'inverse généralisé de  $\hat{\Gamma}_{q,\ell}$  :  $\hat{\rho}_{q,\ell} = \hat{\Gamma}_{q,\ell}^{-1} \hat{\Delta}_{q,\ell}$ .
- Prévision de  $\mathbf{z}_{n+1}$  :  $\hat{\mathbf{z}}_{n+1} = \hat{\rho}_{q,\ell} \hat{\mathbf{x}}_n + \mathbf{A} \ell \bar{\mathbf{z}}$ .



### 3.3.3 Prédiction d'El Nino







Comparaison de l'année 1986 de El Niño avec différentes prévisions.

*Comparaisons des erreurs quadratiques et erreurs absolues relatives de prévision sur la période 1987-1996.*

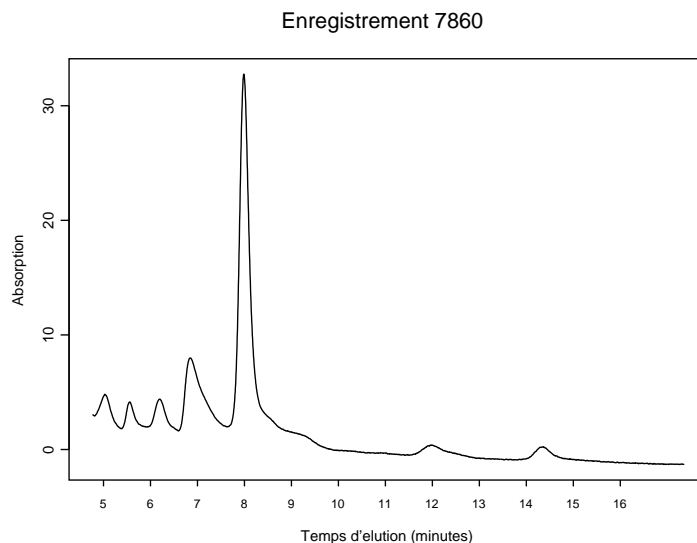
Prédicteur	El Niño index		S. Osc. index	
	EQMP	EAMR	EQMP	EAMR
Climat.	0.73	2.5 %	0.91	6.3 %
SARIMA	1.45	3.7 %	0.95	6.2 %
Noyau	0.60	2.3%	0.87	6.1%
Noyau fonct.	0.58	2.2%	0.82	6.0%
ARH(1) lisse	0.55	2.3%	<b>0.78</b>	<b>5.8%</b>
ARH(1) local	<b>0.53</b>	<b>2.2%</b>	0.82	5.8%

## 3.4 Modélisation par arbre

### 3.4.1 Objectif

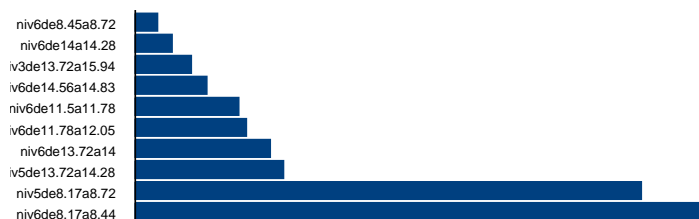
Besse et Farkas (2003)

- Discriminer des jus d'orange présence/absence de pasteurisation.



- Quelle **base** : Fourier, fonctions propres, splines ou ondelettes ?
  - Quelles méthode de **discrimination** : analyse discriminante, régression logistique,  $k$  plus proches voisins, réseaux de neurones, arbres binaires (CART)...
  - **Dimensionnalité** : 6912 mesures pour 64 courbes
  - Problème de **localisation** du **pic** : ondelettes.
  - Comparaison des méthodes de **Sélection** :
  - Procédure de comparaison ( $B$  fois)
1. extraction aléatoire d'un **échantillon test**,
  2. **estimation** du modèle sur l'autre partie (**apprentissage**),
  3. **optimisation** du modèle par validation croisée (sélection de variable ou élagage),
  4. estimation de l'**erreur** sur l'**échantillon test**.

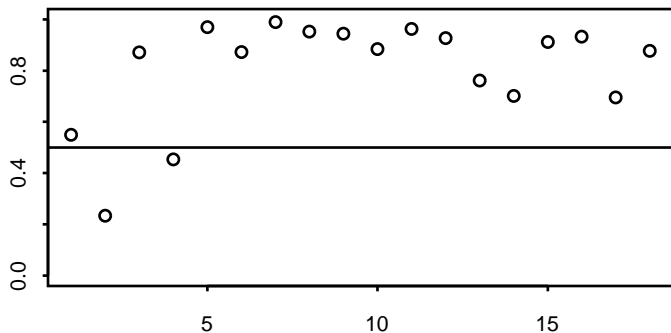
- Échec de la régression logistique.
- Grande instabilité d'un arbre binaire seul.
- Forêt aléatoire avec graphe d'importance :



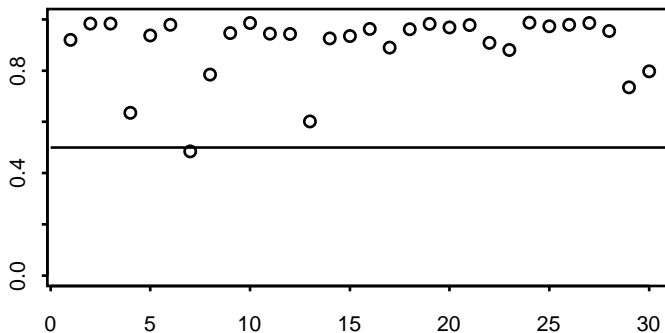
*Importance des variables.*

### jeuVC n° 2 , mtry= 19 , ntrees= 1000

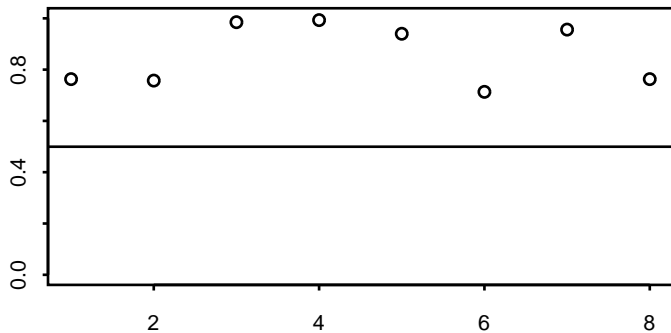
u.s. Train NonPast.: 11.1/100 mc



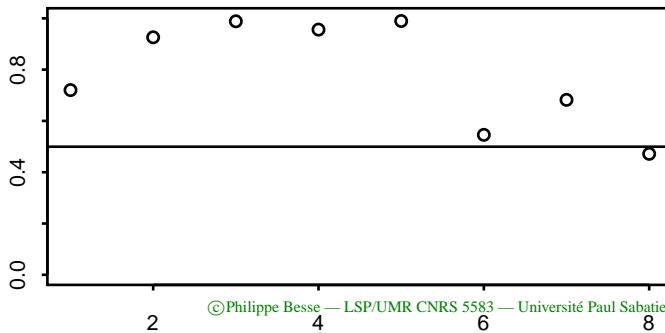
u.s. Train Pasteurises: 3.3/100 mc



u.s. Test NonPast.: 0/100 mc



u.s. Test Pasteurises: 12.5/100 mc



## 3.5 Conclusion

- Coût/intérêt de la démarche.
1. Plongement dans un **espace fonctionnel** (splines, Fourier, ondelettes) adaptée au problème posé : hypothèse ou non de **régularité** par opposition à la présence jugée significative de **singularités**.
  2. **Recaler** les courbes par une transformation non linéaire de l'échelle des temps (*curve registration*).
  3. **Dimensionnalité** et choix de la technique de modélisation ; régularisation (lissage), ACP ou sélection.
  4. **Optimisation** conjointe des paramètres.

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Exploration</b>	<b>5</b>
2.1	ACP de courbes bruitées . . . . .	5
2.1.1	Données . . . . .	5
2.1.2	Modèle et estimation . . . . .	6
2.2	Exemples : ACP de séries climatiques . . . . .	9
2.2.1	ACP des précipitations . . . . .	9
2.2.2	ACP de températures . . . . .	12
<b>3</b>	<b>Modélisation de données fonctionnelles</b>	<b>14</b>
3.1	Introduction . . . . .	14
3.2	Modèle linéaire et extensions . . . . .	15
3.2.1	Modèle linéaire fonctionnel . . . . .	16
3.2.2	Régression sur composantes principales . . . . .	17
3.2.3	Approche par pénalisation . . . . .	18
3.2.4	Modèle linéaire général et autres extensions . . . . .	19
3.3	Prévision fonctionnelle . . . . .	20
3.3.1	Problème . . . . .	20
3.3.2	Modèle ARH(1) . . . . .	21
3.3.3	Prévision d'El Nino . . . . .	24
3.4	Modélisation par arbre . . . . .	27



---

3.4.1	Objectif . . . . .	27
3.5	Conclusion . . . . .	30