

# Régression spline & GAM

Vincent Lefieux



Splines

Interpolation  
spline

Régression spline

GAM

Références

# Plan

Splines

Interpolation spline

Régression spline

GAM

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Plan

## Splines

### Splines

Interpolation  
spline

Régression spline

GAM

Références

# Point de vue physique I

- ▶ Le terme *spline* (cerce en français) désigne une latte en bois flexible utilisée par les dessinateurs industriels pour matérialiser des lignes à courbure variable passant par des points fixés à priori.
- ▶ L'enjeu était d'obtenir des courbes « lisses », d'éventuelles discontinuités pouvant être synonymes de ruptures potentielles à cause d'une faiblesse mécanique.
- ▶ Le tracé de la spline minimise l'énergie de déformation de la latte considérée.

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Point de vue physique II

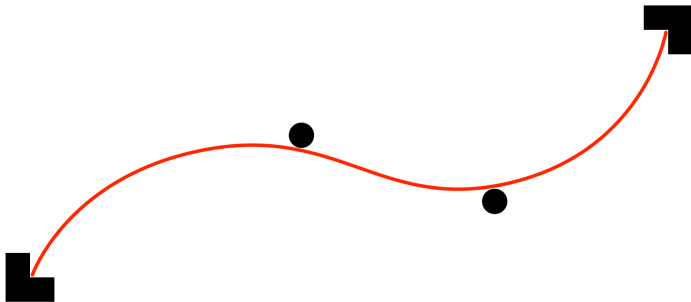
## Splines

Interpolation  
spline

Régression spline

GAM

Références



- Une **spline** est une **fonction définie par morceaux par des polynômes**.

# Interpolation : exemple I



## Splines

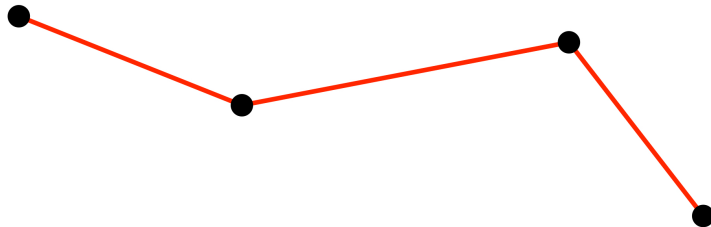
Interpolation  
spline

Régression spline

GAM

Références

# Interpolation : exemple II



## Splines

Interpolation  
spline

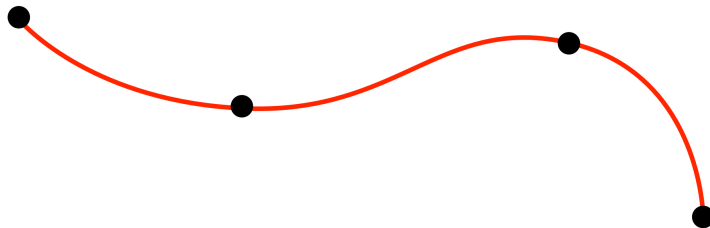
Régression spline

GAM

Références



# Interpolation : exemple III



## Splines

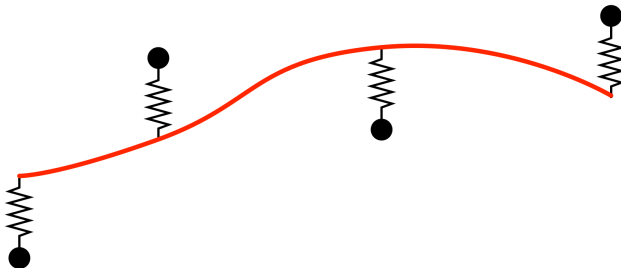
Interpolation  
spline

Régression spline

GAM

Références

# Lissage : point de vue physique I



## Splines

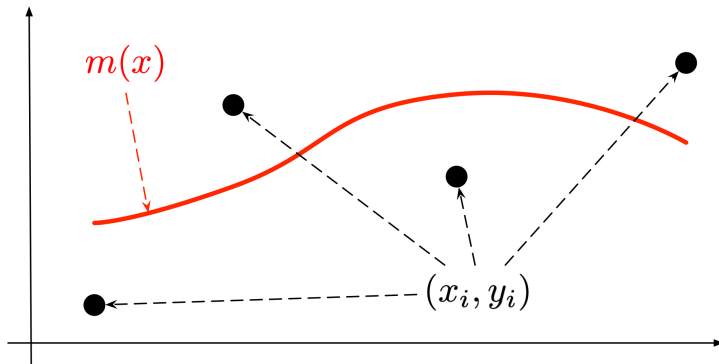
Interpolation  
spline

Régression spline

GAM

Références

# Lissage : point de vue physique II



Splines

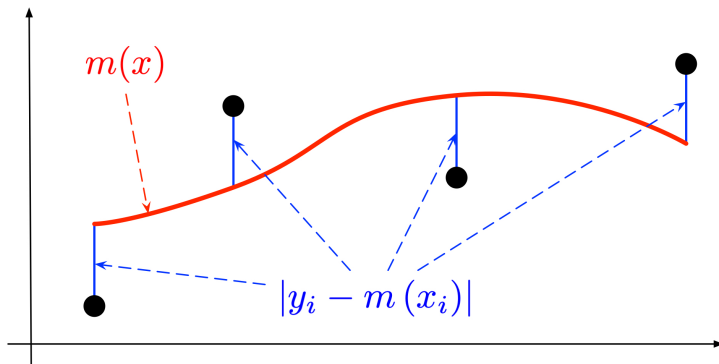
Interpolation  
spline

Régression spline

GAM

Références

# Lissage : point de vue physique III



## Splines

Interpolation  
spline

Régression spline

GAM

Références

## Lissage : point de vue physique IV

- ▶ La **latte**, liée aux points via des ressorts, correspond à une courbe lissée.
- ▶ Sa forme est celle qui **minimise l'énergie totale de déformation** (allongement des ressorts et courbure de la latte) :

$$E(m, \lambda) = \sum_{i=1}^n (y_i - m(x_i))^2 + \lambda \int [m^{(2)}(x)]^2 dx$$

où  $(x_i, y_i)_{i \in \{1, \dots, n\}}$  est l'ensemble des points,  $m$  l'équation de la forme de la latte (avec  $m^{(2)}$  comme dérivée seconde) et  $\lambda$  le rapport entre la raideur de la latte et celle des ressorts.

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Données considérées dans le cas univarié

Splines

Interpolation  
spline

Régression spline

GAM

Références

- ▶ On dispose d'un échantillon de  $(X, Y)$  :

$$\mathcal{D}_n = (X_i, Y_i)_{i \in \{1, \dots, n\}}$$

où  $X \in \mathbb{R}$  et  $Y \in \mathbb{R}$ .

- ▶ On note :

$$d_n = (x_i, y_i)_{i \in \{1, \dots, n\}} \cdot$$

# Splines d'interpolation, de moindres carrés et de lissage

- Pour l'**interpolation**, on utilise des **splines d'interpolation** pour ajuster  $m$  telle que :

$$\forall i \in \{1, \dots, n\} : y_i = m(x_i) .$$

- Pour le modèle de **régression** :

$$\forall i \in \{1, \dots, n\} : y_i = m(x_i) + \varepsilon_i ,$$

on utilise :

- Des **splines de moindres carrés** qui minimisent :

$$\sum_{i=1}^n (y_i - m(x_i))^2 .$$

- Des **splines de lissage** qui minimisent :

$$\sum_{i=1}^n (y_i - m(x_i))^2 + \lambda \int [m^{(2)}(x)]^2 dx .$$

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Données considérées dans le cas multivarié

Splines

Interpolation  
spline

Régression spline

GAM

Références

- On dispose d'un échantillon de  $(X, Y)$  :

$$\mathcal{D}_n = (X_i, Y_i)_{i \in \{1, \dots, n\}}$$

où  $X = (X^1, \dots, X^p)^\top \in \mathbb{R}^p$  et  $Y \in \mathbb{R}$ .

- On note :

$$d_n = (x_i, y_i)_{i \in \{1, \dots, n\}} \cdot$$



# Cas de la régression multivariée

## Splines

Interpolation  
spline

Régression spline

GAM

Références

- ▶ Le modèle **MARS** (Multivariate Adaptive Regression Splines) : (Friedman, 1991).
- ▶ Le modèle **GAM** (Generalized Additive Models) : (Hastie et Tibshirani, 1986).

# Splines polynomiales d'ordre $d$

- ▶ On considère  $K$  **nœuds** (*knots*) :  $(\xi_1, \dots, \xi_K)$  sur  $[a, b]$  :

$$a < \xi_1 < \dots < \xi_K < b .$$

- ▶ Une **spline polynomiale d'ordre  $d \in \mathbb{N}^*$**  est une fonction :

- ▶ **continûment différentiable** jusqu'à l'ordre :

- ▶  $d - 2$  si  $d > 1$ ,
- ▶ 0 si  $d = 1$  (simplement continue),

- ▶ constituée de **polynômes de degré** (inférieur ou égal à)  $(d - 1)$  sur les intervalles inter-nœuds  $[a, \xi_1]$ ,  $[\xi_1, \xi_2], \dots, [\xi_{K-1}, \xi_K], [\xi_K, b]$ .

- ▶ On parle de :

- ▶ **spline linéaire** si  $d = 2$ ,
- ▶ **spline cubique** si  $d = 4$ .

Splines

Interpolation  
spline

Régression spline

GAM

Références

- ▶ On considère  $K$  **nœuds** (*knots*) :  $(\xi_1, \dots, \xi_K)$  sur  $[a, b]$  :

$$a < \xi_1 < \dots < \xi_K < b .$$

- ▶ Une **spline linéaire** (spline d'ordre 2) est une fonction :
  - ▶ **continue**,
  - ▶ constituée de **droites** sur les intervalles inter-nœuds  $[a, \xi_1], [\xi_1, \xi_2], \dots, [\xi_{K-1}, \xi_K], [\xi_K, b]$ .

- ▶ On considère  $K$  **nœuds** (*knots*) :  $(\xi_1, \dots, \xi_K)$  sur  $[a, b]$  :

$$a < \xi_1 < \dots < \xi_K < b .$$

- ▶ Une **spline cubique** (spline d'ordre 4) est une fonction :
  - ▶ **continûment différentiable** jusqu'à l'ordre 2,
  - ▶ constituée de **polynômes de degré** (inférieur ou égal à) **3** sur les intervalles inter-noeuds  $[a, \xi_1]$ ,  $[\xi_1, \xi_2]$ ,  $\dots$ ,  $[\xi_{K-1}, \xi_K]$ ,  $[\xi_K, b]$ .

# Espace des splines polynomiales d'ordre $d$

- ▶ On note  $\mathcal{S}_d(\xi_1, \dots, \xi_K)$  l'ensemble des splines polynomiales d'ordre  $d$  ayant pour nœuds  $(\xi_1, \dots, \xi_K)$ .
- ▶  $\mathcal{S}_d(\xi_1, \dots, \xi_K)$  est un sous-espace vectoriel de l'espace des fonctions dérivables jusqu'à l'ordre  $(d - 2)$  (si  $d > 1$ , 0 sinon), de dimension  $d + K$ .
- ▶ On peut considérer comme base :

$$S_1(x) = 1,$$

$$\vdots$$

$$S_d(x) = x^{d-1},$$

$$\forall k \in \{1, \dots, K\} : S_{d+k}(x) = [(x - \xi_k)_+]^{d-1}$$

où :

$$x_+ = \max(x, 0) = \begin{cases} x & \text{si } x \geq 0 \\ 0 & \text{sinon} \end{cases}.$$

Splines

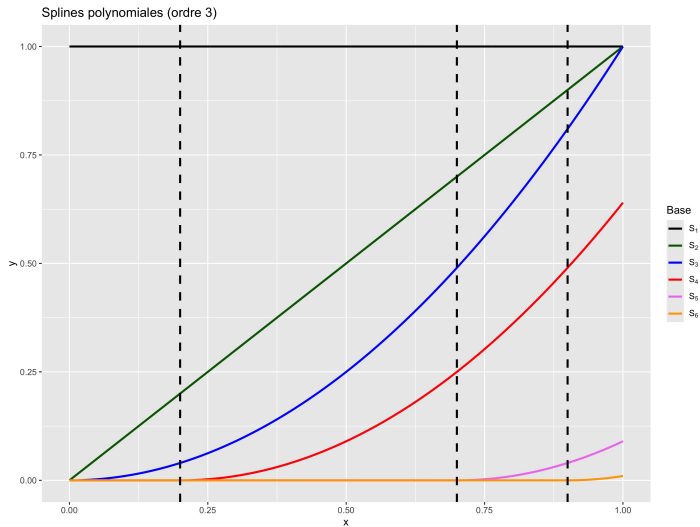
Interpolation  
spline

Régression spline

GAM

Références

# Illustration



## Splines

Interpolation  
spline

Régression spline

GAM

Références

# Remarques

- ▶ On trouve également les **splines naturelles** d'ordre pair : elles diffèrent des splines d'ordre  $d$  au niveau des intervalles  $[a, \xi_1]$  et  $[\xi_K, b]$  sur lesquels elles coïncident avec un polynôme de degré  $\frac{d}{2} - 1$ .  
Des splines cubiques naturelles coïncident donc avec des droites sur le premier et le dernier intervalle.
- ▶ La base  $(S_1, \dots, S_{d+K})$  de  $\mathcal{S}_d(\xi_1, \dots, \xi_K)$ , définie précédemment, est simple d'un point de vue conceptuel mais est peu utilisée en pratique, à cause du support non compact des fonctions de la base, et des problèmes d'arrondis pouvant apparaître pour de grandes valeurs de  $x$ .
- ▶ On lui préfère très souvent la base **B-splines** dans laquelle chaque fonction de base a un support fini. Les B-splines sont une généralisation des courbes de Bézier, et ont été généralisées par les NURBS (Non-Uniform Rational Basis Splines).

Splines

Interpolation  
spline

Régression spline

GAM

Références

# B-splines

- ▶ On considère  $K$  **nœuds** (*knots*) :  $(\xi_1, \dots, \xi_K)$  sur  $[a, b] := [\xi_0, \xi_{K+1}]$  :

$$a = \xi_0 \leq \xi_1 \leq \dots \leq \xi_K \leq b = \xi_{K+1} .$$

Quand les nœuds sont équidistants, on parle de B-splines uniformes.

- ▶ On définit de manière récursive la base de **B-splines** d'ordre  $d \in \mathbb{N}^*$  avec  $d \leq K$  :

1. Pour  $i \in \{1, \dots, K-1\}$  :

$$B_{i,1}(x) = \begin{cases} 1 & \text{si } x \in [\xi_i, \xi_{i+1}[ \\ 0 & \text{sinon} \end{cases} .$$

Si  $\xi_i = \xi_{i+1}$ , on pose par convention  $B_{i,1} = 0$ .

2. Pour  $i \in \{1, \dots, K-d\}$  :

$$B_{i,d}(x) = \frac{x - \xi_i}{\xi_{i+d-1} - \xi_i} B_{i,d-1}(x) + \frac{\xi_{i+d} - x}{\xi_{i+d} - \xi_{i+1}} B_{i+1,d-1}(x) .$$

Par convention, une fraction dont le dénominateur est nul, est considérée nulle.

Splines

Interpolation  
spline

Régression spline

GAM

Références



# Illustration



## Splines

Interpolation  
spline

Régression spline

GAM

Références

# Splines naturelles d'ordre $d$

Splines

Interpolation  
spline

Régression spline

GAM

Références

- ▶ On considère  $K$  **noeuds** (*knots*) :  $(\xi_1, \dots, \xi_K)$  sur  $[a, b]$  :

$$a < \xi_1 < \dots < \xi_K < b .$$

- ▶ Une **spline naturelle d'ordre  $d \in \mathbb{N}^*$  pair** est une fonction :
  - ▶ **continûment différentiable** jusqu'à l'ordre  $d - 2$  ( $d \geq 2$ ),
  - ▶ constituée de **polynômes de degré** :
    - ▶  $\frac{d}{2} - 1$  sur les intervalles  $[a, \xi_1]$  et  $[\xi_K, b]$ ,
    - ▶ (inférieur ou égal à)  **$(d - 1)$**  sur les intervalles inter-noeuds  $[\xi_1, \xi_2], \dots, [\xi_{K-1}, \xi_K]$ .
- ▶ On parle de **spline cubique naturelle** si  $d = 4$ .

# Splines cubiques naturelles

## Splines

### Interpolation spline

### Régression spline

### GAM

### Références

- ▶ On considère  $K$  **noeuds** (*knots*) :  $(\xi_1, \dots, \xi_K)$  sur  $[a, b]$  :

$$a < \xi_1 < \dots < \xi_K < b .$$

- ▶ Une **spline cubique naturelle** est une fonction :
  - ▶ **continûment différentiable** jusqu'à l'ordre 2,
  - ▶ constituée de **polynômes de degré** :
    - ▶ **1** sur les intervalles  $[a, \xi_1]$  et  $[\xi_K, b]$ ,
    - ▶ (inférieur ou égal à) **3** sur les intervalles inter-noeuds  $[\xi_1, \xi_2], \dots, [\xi_{K-1}, \xi_K]$ .

# Espace des splines naturelles d'ordre $d$

## Splines

Interpolation  
spline

Régression spline

GAM

Références

- ▶ On note  $\mathcal{S}_d^*(\xi_1, \dots, \xi_K)$  l'ensemble des splines naturelles d'ordre  $d$  ayant pour nœuds  $(\xi_1, \dots, \xi_K)$ .
- ▶  $\mathcal{S}_d(\xi_1, \dots, \xi_K)$  est un sous-espace vectoriel de l'espace des fonctions dérivables jusqu'à l'ordre  $(d - 2)$ , de dimension  $K$ .

# Plan

## Interpolation spline

Splines

**Interpolation  
spline**

Régression spline

GAM

Références

# Concepts généraux

- ▶ On dit qu'une fonction  $f$  est **absolument continue** s'il existe  $a \in \mathbb{R}$  et une fonction  $g$  intégrable tels que :

$$\forall x \in \mathbb{R} : f(x) = \int_a^x g(y) dy .$$

- ▶ Soit  $W^d(a, b)$  l'ensemble des fonctions  $f$  définies sur  $[a, b]$  telles que :
  - ▶  $(f^{(0)}, \dots, f^{(d-1)})$  sont absolument continues et de carré intégrable,
  - ▶  $f^{(d)}$  est de carré intégrable.
- ▶ La régularité d'une fonction  $f \in W^d(a, b)$  peut être mesurée par :

$$\int_a^b \left[ f^{(d)}(x) \right]^2 dx .$$

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Splines d'interpolation I

Splines

Interpolation  
spline

Régression spline

GAM

Références

- On dispose d'un échantillon de  $(X, Y) \in \mathbb{R} \times \mathbb{R}$  :

$$d_n = (x_i, y_i)_{i \in \{1, \dots, n\}}$$

tel que les  $(x_i)_{i \in \{1, \dots, n\}}$  sont distincts.

- Il existe une unique fonction  $\hat{m} \in W^d(a, b)$  vérifiant :

$$\forall i \in \{1, \dots, n\} : y_i = \hat{m}(x_i) ,$$

$$\hat{m} = \arg \min_{m \in W^d(a, b)} \int_a^b \left[ m^{(d)}(x) \right]^2 dx .$$

- La fonction  $\hat{m}$  est une spline naturelle d'ordre  $2d$  ayant pour nœuds  $(x_1, \dots, x_n)$ .

## Splines d'interpolation II

- Il est possible de relaxer les conditions d'interpolation, notamment dans le cas où les  $(x_i)_{i \in \{1, \dots, n\}}$  ne sont pas distincts, par :

$$\sum_{i=1}^n (y_i - m(x_i))^2 \leq \varepsilon$$

où  $\varepsilon \in \mathbb{R}^*$  donne le niveau de la relaxation.

- On peut montrer que la fonction  $\hat{m}$  est alors la solution du problème d'optimisation suivant :

$$\hat{m} = \arg \min_{m \in W^d(a,b)} \sum_{i=1}^n (y_i - m(x_i))^2 + \lambda \int_a^b [m^{(d)}(x)]^2 dx$$

où  $\lambda$  dépend de  $\varepsilon$  (de manière complexe).

Splines

Interpolation  
spline

Régression spline

GAM

Références



# Plan

Régression spline

Splines

Interpolation  
spline

**Régression spline**

GAM

Références

- ▶ On dispose d'un échantillon de  $(X, Y)$  :

$$\mathcal{D}_n = (X_i, Y_i)_{i \in \{1, \dots, n\}}$$

où  $X \in \mathbb{R}$  et  $Y \in \mathbb{R}$ .

- ▶ On note :

$$d_n = (x_i, y_i)_{i \in \{1, \dots, n\}} \cdot$$

- On appelle **spline de moindres carrés d'ordre  $d$**  ayant comme nœuds  $(\xi_1, \dots, \xi_K)$ , la fonction  $\hat{m}$  suivante :

$$\hat{m} = \arg \min_{m \in \mathcal{S}_d(\xi_1, \dots, \xi_K)} \sum_{i=1}^n (y_i - m(x_i))^2 .$$

## Splines de moindres carrés II

- ▶ On considère la **base**  $(S_1, \dots, S_{d+K})$  de  $\mathcal{S}_d(\xi_1, \dots, \xi_K)$ .
- ▶ On peut écrire :

$$\hat{m}(x) = \sum_{j=1}^{d+K} \hat{\theta}_j S_j(x)$$

où  $(\hat{\theta}_1, \dots, \hat{\theta}_{d+K})^\top \in \mathbb{R}^{d+K}$  **minimisent** :

$$\sum_{i=1}^n \left( y_i - \sum_{j=1}^{d+K} \theta_j S_j(x_i) \right)^2 .$$

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Splines de moindres carrés III

- On note :

$$\mathbf{y} = (y_1, \dots, y_n)^\top,$$

$$\hat{\mathbf{y}} = (\hat{m}(x_1), \dots, \hat{m}(x_n))^\top,$$

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_{d+K})^\top,$$

$$\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_{d+K})^\top,$$

$$\mathbf{N} = [S_j(x_i)]_{i \in \{1, \dots, n\}, j \in \{1, \dots, d+K\}}.$$

- On cherche  $\hat{\boldsymbol{\theta}}$  qui minimise :

$$(\mathbf{y} - \mathbf{N}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{N}\boldsymbol{\theta}).$$

- La solution est :

$$\hat{\boldsymbol{\theta}} = (\mathbf{N}^\top \mathbf{N})^{-1} \mathbf{N}^\top \mathbf{y}.$$

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Splines de moindres carrés IV

Splines

Interpolation  
spline

Régression spline

GAM

Références

- On peut écrire :

$$\hat{\mathbf{y}} = \mathbf{S} \mathbf{y}$$

où :

$$\mathbf{S} = \mathbf{N} \left( \mathbf{N}^\top \mathbf{N} \right)^{-1} \mathbf{N}^\top$$

est la **matrice de lissage** (*smoothing*).

# Splines de moindres carrés V

- ▶ Si on considère une base de B-splines, la matrice  $N^T N$  est alors une matrice bande avec  $2d - 1$  diagonales non nulles.
- ▶ On considère **usuellement** :
  - ▶ un **ordre  $d = 4$** ,
  - ▶ un **nombre de noeuds  $K \in \{0, \dots, n - d\}$** ,
  - ▶ des **noeuds  $(\xi_1, \dots, \xi_K)$  équirépartis ou égaux à des quantiles empiriques de  $X$** .
- ▶ A  $d$  fixé, le nombre de noeuds  **$K$  est un hyperparamètre modulant le lissage** :
  - ▶  $K = 0$  : l'estimateur correspond à une régression polynomiale de degré  $d - 1$ .
  - ▶  $K = n - d$  : l'estimateur correspond à une spline d'interpolation.

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Des splines de moindres carrés aux splines de lissage

Splines

Interpolation  
spline

Régression spline

GAM

Références

- ▶ On ajoute une **pénalité** afin de **contrôler les variations de l'estimateur** (importantes dans le cas des splines de moindres carrés).



- On appelle **spline de lissage d'ordre  $2d$**  ayant comme nœuds  $(\xi_1, \dots, \xi_K)$ , la fonction  $\hat{m}$  suivante :

$$\hat{m} = \arg \min_{m \in W^d(a,b)} \sum_{i=1}^n (y_i - m(x_i))^2 + \lambda \int_a^b \left[ m^{(d)}(x) \right]^2 dx$$

où  $\lambda \in \mathbb{R}^{+*}$  caractérise le compromis entre l'ajustement et le caractère lisse de la fonction.

- Si on considère la **base**  $(S_1, \dots, S_n)$  de  $\mathcal{S}_{2d}^*(x_1, \dots, x_n)$  (splines naturelles), on peut alors écrire :

$$\hat{m}(x) = \sum_{j=1}^n \hat{\theta}_j S_j(x)$$

où  $(\hat{\theta}_1, \dots, \hat{\theta}_n) \in \mathbb{R}^n$  minimisent :

$$\sum_{i=1}^n \left( y_i - \sum_{j=1}^n \theta_j S_j(x_i) \right)^2 + \lambda \int_a^b \left( \sum_{j=1}^n \theta_j S_j^{(d)}(x) \right)^2 dx.$$

# Splines de lissage III

- On note :

$$\mathbf{y} = (y_1, \dots, y_n)^\top ,$$

$$\hat{\mathbf{y}} = (\hat{m}(x_1), \dots, \hat{m}(x_n))^\top ,$$

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)^\top ,$$

$$\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_n)^\top ,$$

$$\mathbf{N} = [S_j(x_i)]_{i \in \{1, \dots, n\}, j \in \{1, \dots, n\}} ,$$

$$\boldsymbol{\Omega} = \left[ \int_a^b S_i^{(d)}(x) S_j^{(d)}(x) dx \right]_{i \in \{1, \dots, n\}, j \in \{1, \dots, n\}} .$$

- On cherche  $\hat{\boldsymbol{\theta}}$  qui minimise :

$$(\mathbf{y} - \mathbf{N}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{N}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^\top \boldsymbol{\Omega} \boldsymbol{\theta} .$$

- La solution est :

$$\hat{\boldsymbol{\theta}} = (\mathbf{N}^\top \mathbf{N} + \lambda \boldsymbol{\Omega})^{-1} \mathbf{N}^\top \mathbf{y} .$$

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Splines de lissage IV

Splines

Interpolation  
spline

Régression spline

GAM

Références

- On peut écrire :

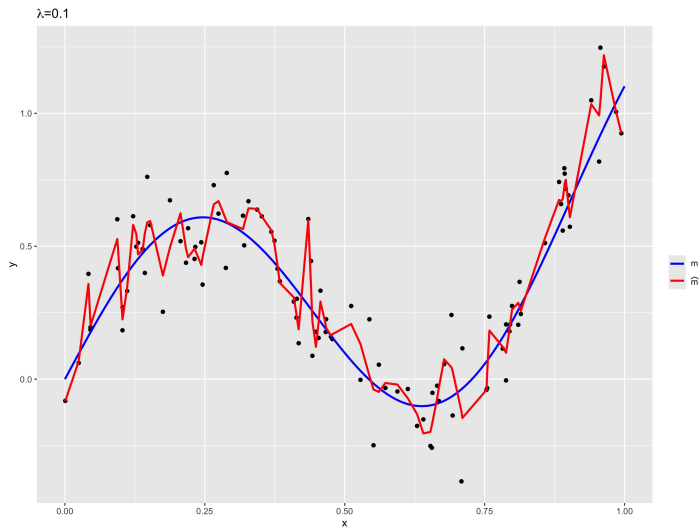
$$\hat{\mathbf{y}} = \mathbf{S} \mathbf{y}$$

où :

$$\mathbf{S} = \mathbf{N} \left( \mathbf{N}^{\top} \mathbf{N} + \lambda \mathbf{\Omega} \right)^{-1} \mathbf{N}^{\top}$$

est la **matrice de lissage** (*smoothing*).

# Illustration I



Splines

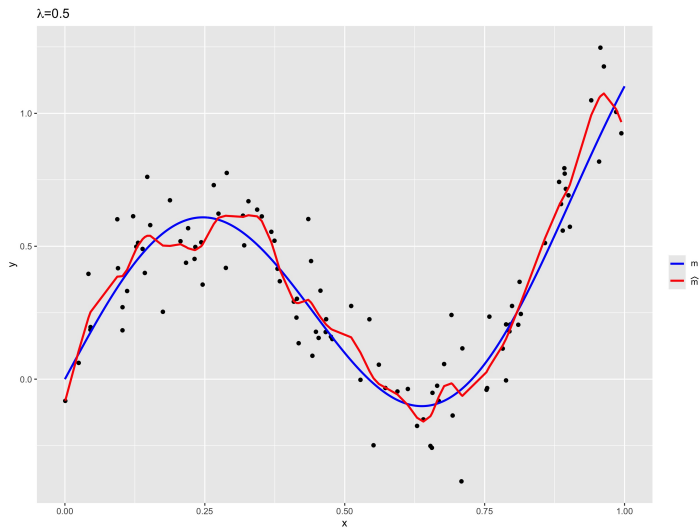
Interpolation  
spline

Régression spline

GAM

Références

# Illustration II



Splines

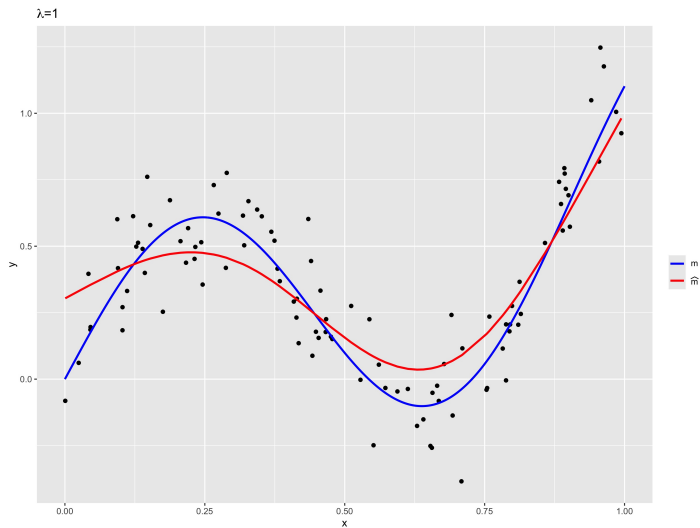
Interpolation  
spline

Régression spline

GAM

Références

# Illustration III



Splines

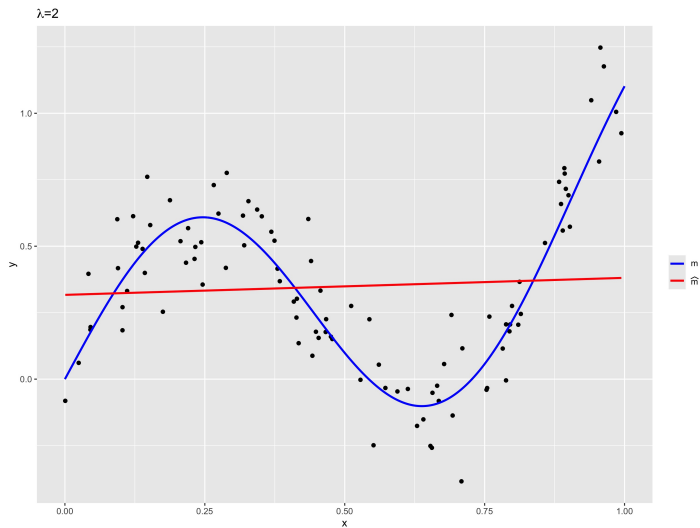
Interpolation  
spline

Régression spline

GAM

Références

# Illustration IV



Splines

Interpolation  
spline

Régression spline

GAM

Références



# Plan

## GAM

Splines

Interpolation  
spline

Régression spline

**GAM**

Références

- ▶ On dispose d'un échantillon de  $(X, Y)$  :

$$\mathcal{D}_n = (X_i, Y_i)_{i \in \{1, \dots, n\}}$$

où  $X = (X^1, \dots, X^p)^\top \in \mathbb{R}^p$  et  $Y \in \mathbb{R}$ .

- ▶ On note :

$$d_n = (x_i, y_i)_{i \in \{1, \dots, n\}} \cdot$$

# Le modèle

- ▶ Le **modèle additif généralisé** (**GAM** : *Generalized Additive Model*) suppose que :

$$Y = c + \sum_{j=1}^p g_j(X^j) + \varepsilon$$

où  $g_j : \mathbb{R} \rightarrow \mathbb{R}$  sont  $p$  fonctions inconnues.

- ▶ Pour assurer l'existence, il faut imposer une contrainte, par exemple :

$$\forall j \in \{1, \dots, p\} : \int g_j(x) dx = 0 .$$

- ▶ On peut notamment estimer les fonctions  $(g_j)_{j \in \{1, \dots, p\}}$  par noyau, par polynômes locaux, par projection sur des bases orthogonales et par ajustement spline.
- ▶ La méthode des **splines** est la plus couramment employée pour les modèles GAM.

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Le critère d'estimation

- ▶ On suppose ici que chaque fonction  $g_j$  est estimée à l'aide de fonctions splines.
- ▶ On cherche à minimiser :

$$\sum_{i=1}^n \left( y_i - c - \sum_{j=1}^p g_j(x_i^j) \right)^2 + \sum_{j=1}^p \lambda_j \int \left[ g_j^{(2)}(x) \right]^2 dx$$

où  $(\lambda_j)_{j \in \{1, \dots, p\}} \in \mathbb{R}^{+p}$  sont des hyperparamètres de régularisation.

- ▶ Les solutions sont des splines cubiques, chaque fonction  $g_j$  ayant pour noeuds les  $(x_i^j)_{i \in \{1, \dots, n\}}$ .
- ▶ On impose comme contrainte d'unicité :

$$\forall j \in \{1, \dots, p\} : \sum_{i=1}^n g_j(x_i^j) = 0 .$$

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Estimation : méthode de backfitting

## 1. Initialisation :

$$\hat{c} = \bar{y} ,$$

$$\forall j \in \{1, \dots, p\} : \hat{g}_j(x) = 0 .$$

## 2. Pour $k \in \{1, \dots, p\}$ :

2.1 On **estime**  $g_k$  en fixant tous les autres (étape de **backfitting**).

Le problème à minimiser est :

$$\min_{g_k} \sum_{i=1}^n \left( y_i - \hat{c} - \sum_{j=1, j \neq k}^p \hat{g}_j(x_i^j) - g_k(x_i^k) \right)^2 + \lambda_k \int [g_k^{(2)}(x)]^2 dx .$$

2.2 On **centre**  $\hat{g}_k$  en lui soustrayant :

$$\frac{1}{n} \sum_{i=1}^n \hat{g}_k(x_i^k) .$$

On **itère l'étape 2** jusqu'à stabilisation de l'optimisation.

Splines

Interpolation  
spline

Régression spline

GAM

Références

- ▶ Il existe d'autres méthodes d'estimation pour les modèles GAM.
- ▶ On peut réunir certaines covariables (en veillant au fléau de la dimension).
- ▶ On peut utiliser les modèles sur des **séries temporelles**.

Splines

Interpolation  
spline

Régression spline

GAM

Références

# Mise en œuvre logicielle

- ▶ Sous **Python** :
  - ▶ Package **pygam**.
- ▶ Sous **R** :
  - ▶ Package **gam** proposé par Hastie.
  - ▶ Package **mgcv** proposé par Wood.

Splines

Interpolation  
spline

Régression spline

**GAM**

Références

Friedman, J. H. 1991, «Multivariate adaptive regression splines», *The Annals of Statistics*, vol. 19, n° 1, p. 1–67.

Hastie, T. et R. Tibshirani. 1986, «Generalized additive models», *Statistical Science*, vol. 1, n° 3, p. 295–318.