

## Contents

<b>1</b>	<b>Introduction générale</b>	<b>2</b>
1.1	Motivation . . . . .	2
1.2	Problématique . . . . .	2
<b>2</b>	<b>Méthodes non-paramétriques</b>	<b>3</b>
2.1	Définitions . . . . .	3
2.2	Estimateurs par projection : . . . . .	3
<b>3</b>	<b>Estimateur de densité à noyau :</b>	<b>4</b>
3.1	Évaluer un estimateur . . . . .	4
3.2	Méthodes adaptatives : . . . . .	5
3.2.1	Choix du noyau . . . . .	5
3.2.2	Choix de la fenêtre . . . . .	5

---

# 1 Introduction générale

## 1.1 Motivation

On a une *arbre phylogénétique* présentant les relations de parenté entre *espèces* et on s'intéresse au *branchement évolutif* apparaît après une durée aléatoire d'une loi fixée  $\mu$  indépendamment du passé et du futur évolutifs des espèces .

Quelle est cette loi  $\mu$  ? Sa variance ? Sa moyenne ?.

## 1.2 Problématique

On observe les successions de branchements qui composent un arbre phylogénétique et à partir de ces données quantitatives observées on veut estimer la fonction  $f$  qui donne le temps qu'il faut pour qu'un branchement évolutif apparaisse .

Formellement ; On a un échantillon  $\chi = \{X_1, \dots, X_n\}$  de variables observées qui ont pour fonction de densité  $f$  dans  $\mathcal{F}$  où  $\mathcal{F}$  est un espace fonctionnel . On cherche à estimer cette fonction de densité  $f$  sur la quelle on fait le moins d'hypothèses possibles . D'où le modèle suivant  $\{P = P_f, f \in \mathcal{F}\}$  qui revient à faire une estimation non-paramétrique .

Pour répondre à la problématique on cherchera à comprendre et finalement implémenter un estimateur à noyau adaptatif , de type Goldenshluger-Lepski sur des données d'arbres phylogénétiques .

## 2 Méthodes non-paramétriques

En statistique, on parle de l'estimation quand cherche à trouver certains paramètres inconnus caractérisant une distribution à partir d'un échantillon de données observées en se basant sur différentes méthodes. On se tourne vers l'estimation non-paramétrique lorsqu'on traite des paramètres à dimension infinie. C'est qui est notre cas, comme on cherche à estimer une fonction densité qui appartient à un espace fonctionnel.

On présente dans la suite une courte introduction de l'estimation non paramétrique. On introduit ensuite les deux classes principales de l'estimation fonctionnelle (l'estimation par projection et l'estimation à noyau) afin de discuter de ces deux classes et de pourquoi on fait le choix de l'estimation à noyau.

### 2.1 Définitions

**Définition 1.** *Estimation non-paramétrique:*

L'estimation non-paramétrique vise à résoudre des problèmes d'estimation dans le cadre statistique où le modèle dont on s'intéresse n'est pas décrit par un nombre fini de paramètres et dont chacun de ces paramètres ne permet pas de décrire la structure générale de la distribution des variables aléatoires.

Cela signifie qu'on utilise des modèles statistiques à dimension infinie.

Dans le cadre de notre problématique on s'intéresse à l'estimation de densité.

**Définition 2.** *Estimation à densité:*

Un des principes de base de l'estimation à densité selon une méthode d'estimation non-paramétrique est le suivant :

Pour un échantillon des observations quantitatives  $\mathbb{X} = X_1, \dots, X_n$  des variables aléatoires i.i.d admettant une densité  $f = F'$ . Supposons que  $f \in \mathcal{F}$  où  $\mathcal{F}$  un espace fonctionnel on cherche à estimer la fonction inconnue  $f$  à partir de ces observations.

On notera  $\hat{f}_n$  l'estimateur de  $f$ .

On se trouve donc avec le modèle suivant  $\mathbb{P} \quad \mathbb{P}_f, f \in \mathcal{F}$ , tel que  $\mathbb{P}_f$  est la mesure probabilité de la densité  $f$ .

L'estimation ici concerne donc la fonction elle-même plutôt que les paramètres, ce qui explique le nom de l'estimation non-paramétrique.

**Remarque 1.** -On notera dans la suite  $\hat{f}$  l'estimateur de la vraie fonction  $f$ .

-On prend souvent la distance  $L^p$  avec  $p = 1, 2$  ou  $\infty$ .

Ils existent deux grandes familles de méthodes pour estimer une fonction densité :

Estimation par projection.

Estimation par le noyau.

### 2.2 Estimateurs par projection :

**Définition 3.** *Estimation par projection:*

Supposant que la fonction  $f$  à estimer est dans l'espace de Hilbert  $\mathcal{F} = (L^2, \|\cdot\|, \langle \cdot, \cdot \rangle)$  avec  $(\Phi_j)_{j \geq 0}$  une base orthonormée de  $L^2$ ,  $\mathbb{E}_N$  un sous-espace fini de  $\mathcal{F}$  et  $1 \leq |N| < \infty$ .

De plus  $a_\lambda = \langle f, \Phi_\lambda \rangle = \int_{\mathbb{R}} f(x) \Phi_\lambda(x) dx$ .

Alors, on estime la fonction  $f$  par son projeté

$$\Pi_N f = \sum_{\lambda \in N} a_\lambda \Phi_\lambda$$

**Remarque 2.** - Cette méthode nous ramène au cas paramétrique .  
 - Plus la valeur de  $N$  est grande plus le biais est petit et la variance est grande .

Dans la suite on procédera avec la méthode la plus fréquente utilisée pour l'estimation d'une densité :  
 L'estimation par le noyau.

### 3 Estimateur de densité à noyau :

#### 3.1 Évaluer un estimateur

Pour évaluer un estimateur on définit le risque associé d'un estimateur  $\hat{f}$  pour l'estimateur  $f$ .

**Définition 4.** La fonction de risque :

$$\mathcal{R}(\hat{f}, f) = \mathbb{E}_f[||\hat{f} - f||^2]$$

**Remarque 3.** La fonction de risque associé nous permet de comparer l'estimateur  $\hat{f}$  et l'estimation  $f$ .  
 On cherche à ce que ce risque associé soit minimal (i.e tend vers 0 pour un nombre d'observation assez grand.)

**Proposition.** Dans le cas d'un estimateur à noyau , on a :

$$R(\hat{f}, f) = E_f[||\hat{f} - f||^2] = ||f - K_h * f||^2 + E_f[||\hat{f} - K_h * f||^2]$$

**Démonstration.** On que :

$$E_f[||\hat{f} - f||^2] = \mathbb{E}_f[||\hat{f} + \mathbb{E}_f(\hat{f}) - (\mathbb{E}_f(\hat{f}) - f)||^2].$$

$$E_f[||\hat{f} - f||^2] = \mathbb{E}_f[||\hat{f} + E_f(\hat{f})||^2] + \mathbb{E}_f[||\mathbb{E}_f(\hat{f}) - f||^2] - 2\mathbb{E}_f(< \hat{f} - \mathbb{E}_f(\hat{f}); \mathbb{E}_f(\hat{f}) - f >).$$

Comme  $\hat{f}$  est déterministe

$$2\mathbb{E}_f(< \hat{f} - E_f(\hat{f}); \mathbb{E}_f(\hat{f}) - f >) = 2 < 0, E_f(\hat{f}) - f > = 0$$

Ainsi  $||\mathbb{E}_f(\hat{f}) - f||$  est déterministe

On obtient :

$$R(\hat{f}, f) = \mathbb{E}_f[||\hat{f} - f||^2] = ||f - K_h * f||^2 + \mathbb{E}_f[||\hat{f} - K_h * f||^2]$$

On à bien trouver que

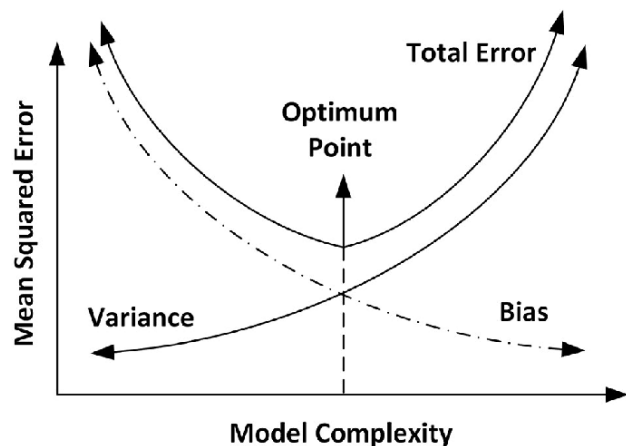
$$R(\hat{f}, f) = \text{biais}^2 + \text{Var}$$

**Remarque 4.** Plus la valeur de  $h$  est grande, plus le biais devient grand et la variance petite et vis-versa ; plus la valeur de  $h$  est petite plus le biais devient petit et la variance s'explode.

Donc afin de minimiser l'expression de risque le choix de  $h$  est très influent et même plus crucial pour la qualité de l'estimateur que celui de la noyau  $K$ .

On doit chercher le meilleur compromis biais-variance pour afin avoir un risque minimal.

Un paramètre trop faible provoque l'apparition de détails artificiels sur le graph de l'estimateur (La variance devient trop grande), par contre si on prend une valeur de  $h$  très grande on aura la majorité des caractéristiques effacée .



### 3.2 Méthodes adaptatives :

On a introduit précédemment la notion de l'estimation de densité qui dépend d'un paramètre de lissage  $h$ . Soit  $(\hat{f}_h)_{h \in \mathcal{H}}$  une famille des estimateurs de la vraie fonction densité  $f$ . La question qui se pose est donc la suivante : comment peut-on construire un estimateur à risque optimal à partir de cette famille (en prenant en considération les observations) ?

Dans cette partie et afin de répondre à la question qu'on a posée on va discuter au premier temps du choix du noyau. Ensuite, on va introduire deux méthodes pour le choix du paramètre de lissage  $h$ .

#### 3.2.1 Choix du noyau

On note par le noyau la fonction intégrable  $K : \mathbb{R} \rightarrow \mathbb{R}$  tel que  $\int_{\mathbb{R}} K(u) du = 1$ .

Soient :

$h > 0$  le paramètre de lissage.

$$K_h : u \in \mathbb{R} \rightarrow K\left(\frac{u}{h}\right)/h.$$

**Lemme 1.** On peut approximer la famille  $(K_h)_{h>0}$  par l'identité du produit de convolution.

**Démonstration.** A Faire

**Corollaire 1.**  $K_h * f : x \rightarrow \int_{\mathbb{R}} K_h(y-x)f(y)dy$  tend vers la fonction  $f$  quand  $h$  tend vers 0. (pour la distance  $L^2$ )

#### 3.2.2 Choix de la fenêtre

L'estimation de densité nécessite le choix de la fenêtre qu'on note  $h$ .

En statistique non-paramétrique, il existe plusieurs méthodes et critères de qualité pour le choix de la fenêtre.

On présente dans la suite deux méthodes :

Méthode de validation croisée.

Méthode de Goldenshluger-Lepski.

##### 3.2.2.1 Méthode par validation croisée

**3.2.2.2 Méthode de Goldenshluger-Lepski** La méthode de Lepski donne principalement des critères pour le choix entre estimateurs à noyau  $(\hat{f}_h)_{h \in \mathcal{H}}$  avec différentes fenêtres qu'on fixe en prenant en considération l'échantillon des observations.

Cette méthode propose de choisir le  $\hat{h}$  qui minimise l'expression suivante :

$$B(h) + V(h)$$

Avec :

$$B(h) = \sup_{h' \in \mathcal{H}} [|\hat{f}_{h'} - \hat{f}_h| - V(h')]$$

Et

$$V(h) = a \frac{\|K_{h'}\|^2}{n}$$

Tel que K est le noyau ,a un paramètre et V(h) est le terme de pénalisation .

On a donc le  $\hat{h}$  est égale à :

$$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} (B(h) + V(h))$$

**Remarque 5.** *Le terme de pénalisation choisi est proportionne a la variance de l'estimateur .*

On s'intéresse dans cette méthode à déterminer le terme de pénalisation minimal V(h) tel que si on le dépasse on n'obtient plus l'équilibre biais-variance .

Dans ce cas , la valeur de  $\hat{h}$  est d'ordre  $\frac{1}{n}$  ,

Le choix optimal de la fenêtre h suivant cette méthode dans ce cas est  $n^{-\frac{1}{2\alpha+1}}$