

Pour commencer

Sophie Manuel Stéphane Sadio Wiam CHAOUI

Contents

Motivation :	2
Remarques/Questions :	2
Statistique non-paramétrique . Pourquoi ? :	2
Estimateur de densité à noyau , Pourquoi ? Pourquoi pas d'autres méthodes ?.	2
Evaluer un estimateur :	2
Méthodes adaptatives :	3
Methode de Goldenshluger-Lepski :	3
P.S:	3

Motivation :

On a une arbre phylogénétique, qu'un branchement évolutif (i.e. la création d'une espèce) apparaît après une durée aléatoire d'une loi fixée μ indépendamment du passé et du futur évolutifs des espèces. On cherche cette loi μ ? sa variance ? sa moyenne ?

(mots clés: Arbre phylogénétique, branchement évolutif, loi fixée)

Remarques/Questions :

C'est quoi le but/la question principale du projet ?

Pourquoi l'estimation non-paramétrique ? (le lien avec la création d'une espèce et les arbres phylogénétiques)

Pourquoi l'estimation de densité à noyau ?

En pratique toute fonction mesurable \hat{f} des data est un estimateur de f , comment donc peut-on juger (évaluer) les performances de ces estimateurs afin de choisir un ? (indication : Risque quadratique $(R(\hat{f}, f))$)

Statistique non-paramétrique . Pourquoi ? :

Dans notre cas on a des données observées quantitatives présenter par les successions de branchements qui composent un arbre phylogénétique .

Soit un arbre phylogénétique on cherche à estimer la fonction f qui donne la durée qui faut pour qu'un branchement évolutif apparaisse. (qui a généré l'échantillon aléatoire)

Discription formelle : $\chi = \{X_1, \dots, X_n\}$ un échantillon de variables observées qui ont pour fonction de densité $f \in F$ où F est un espace fonctionnel et a partie des observations (data) X_1, \dots, X_n on veut estimer cette fonction densité f sur la quelle on fait le moins d'hypothèses possibles . On a alors le modèle suivant $\{P = P_f, f \in F\}$. (Ce qui revient à une estimation non-paramétrique)

Estimateur de densité à noyau , Pourquoi ? Pourquoi pas d'autres méthodes ?.

Un noyau est une fonction intégrable qu'on note $K : R \rightarrow R$ et qui vérifie $\int_R K(u) du = 1$.

Pourquoi l'estimation à noyau est plus intéressante ?

Soient $h > 0$ et $K_h : u \in R \rightarrow \frac{K(\frac{u}{h})}{h}$ alors pour la famille $(K_h)_{h \geq 0}$ on a le résultat de convolution suivant : $K_h * f : x \rightarrow \int_R K_h(x - x') f(x') dx'$ tend vers f quand h tend vers 0 .

La densité f peut alors être estimée par le produit de convolution $K_h * f$ (qui satisfait $K_h * f = E_f[K_h(x - X_1)]$)

Donc l'estimateur à noyau de f pour un $h > 0$ fixé est :

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - X_i) \quad x \in R$$

Evaluer un estimateur :

Pour cela il faut définir le risque associé d'une estimation \hat{f} pour l'estimation de f .

La fonction de risque est :

$$R(\hat{f}, f) = E_f[||\hat{f} - f||^2] \text{ (afin de comparer les deux fonctions)}$$

Faut se poser la question sur le choix de la distance (la norme) ?

On prend souvent la distance L^p pour $p = 1, 2$ ou ∞

On veut que le risque soit minimal, tend vers 0 pour un nombre d'observations assez grand

Dans le cas d'un estimateur à noyau :

$R(\hat{f}, f) = E_f[||\hat{f} - f||^2] = ||f - K_h * f||^2 + E_f[||\hat{f} - K_h * f||^2]$ (pourquoi on a l'égalité ? réponse en P.S à la fin)

(P.10 et 11 de AN INTRODUCTION TO NONPARAMETRIC ADAPTIVE ESTIMATION)

On retrouve $R(\hat{f}, f) \leq ||f - K_h * f||^2 + ||K||^2 \frac{1}{nh}$

Pour minimiser cette dernière expression le choix de h est très influent. (Le choix de h est plus crucial pour la qualité de l'estimateur que celui de K)

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

Méthodes adaptatives :

D'après ce qui précède on a introduit la notion de l'estimation de densité à noyau qui dépend d'un paramètre de lissage h . C'est à dire qu'on a bien défini une famille $(\hat{f}_h)_{h \in \beta_n}$ des estimateurs de la vraie fonction densité f .

Comment peut-on alors construire un estimateur à risque optimal à partir de cette famille $(\hat{f}_h)_{h \in \beta_n}$ en prenant en considération les observations ? (On veut un estimateur adaptatif qui donne le meilleur équilibre biais-variance)

Méthode de Goldenshluger-Lepski :

P.S:

L'estimation non-paramétrique ne fait aucune hypothèse sur la nature/forme/type de la distribution des variables aléatoires/sur l'appartenance de la fonction densité

Plus que le risque on peut juger un estimateur selon son efficacité

On veut que :

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

$$E_f[||\hat{f} - f||^2] = E_f[||\hat{f} + E_f(\hat{f})||^2] + E_f[||E_f(\hat{f}) - f||^2] - 2E_f(\langle \hat{f} - E_f(\hat{f}); E_f(\hat{f}) - f \rangle)$$

Comme \hat{f} est déterministe :

$$2E_f(\langle \hat{f} - E_f(\hat{f}); E_f(\hat{f}) - f \rangle) = 2 \langle \hat{f} - E_f(\hat{f}), E_f(\hat{f}) - f \rangle = 0 \text{ Ainsi } ||E_f(\hat{f}) - f|| \text{ est déterministe :}$$

On obtient :

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