

# Utilisation de méthodes combinatoires en estimation de densité

Groupe de travail 2<sup>ème</sup> année

À partir de l'ouvrage « Combinatorial methods in density estimation »

par Luc Devroye et Gábor Lugosi, 2000

Disponible en 2 exemplaires à la bibliothèque.

En pratique : \* Séances les jeudis après-midis, de 13h30 à 14h45

- exposé durant entre 50 minutes et 1 heure

DÉBUT

LE JEUDI 7  
OCTOBRE !

- suivi de 15 minutes de questions / discussions

- horaire strict car de nombreux jeudis nous avons un séminaire à 15h à l'École des Mines

\* Scannez vos notes manuscrites afin que les suivants aient une trace écrite

---

## Présentation du problème

$X_1 \dots X_n$  observations iid selon une loi  $\mathcal{F}$ , de fonction de répartition  $F$   
réelles

(1) - Comment estimer  $F$  si celle-ci est inconnue ? (problème statistique)

- On note  $\hat{F}_n$  la fonction de répartition empirique, i.e., la fonction de répartition associée à la mesure empirique  $\mathcal{F}_n = \frac{1}{n} \sum_{t=1}^n \delta_{X_t}$  :

$$\forall x \in \mathbb{R}, \quad \hat{F}_n(x) = \frac{1}{n} \sum_{t=1}^n \mathbb{1}_{\{X_t \leq x\}}$$

C'est un bon estimateur : il est consistant, c'est-à-dire qu'il converge en probabilité vers la quantité à estimer  $F$ , pour la distance de la norme du Supremum, et ce, quelle que soit  $F$  :

$$d_{\infty}(\hat{F}_n, F) \stackrel{\text{not.}}{=} \|\hat{F}_n - F\|_{\infty} \xrightarrow{P} 0$$

En fait, la convergence a même lieu p.s. : c'est le théorème de

## Glivenko - Cantelli.

- (2) - La méthode précédente a donc permis d'estimer simultanément et de manière convergente toutes les quantités

$$\nu([-x_0, x]) \quad x \in \mathbb{R}$$

- Que faire si l'on veut désormais faire de même pour tous les
- $$\nu(A), \quad A \in \mathcal{B}(\mathbb{R})$$

où  $\mathcal{B}(\mathbb{R})$  désigne la tribu des boréliens de  $\mathbb{R}$ ?

- (3) - On introduit à cet effet le critère suivant, dit distance en variation :  
soient  $\mu_1$  et  $\mu_2$  deux lois sur  $\mathbb{R}$ ,

(C'est bien une distance.)

$$d_{\text{VT}}(\mu_1, \mu_2) = \sup_{A \in \mathcal{B}(\mathbb{R})} |\mu_1(A) - \mu_2(A)|$$

- Par exemple, lorsque  $\nu$  est diffuse (n'a aucun atome) et désignant toujours par  $\nu_n$  la mesure empirique (qui, elle, est p.s. totalement atomique):

$$d_{\text{VT}}(\nu_n, \nu) = 1.$$

- (4) Simplification et formulation d'une stratégie plus raisonnable

- Simplification: Supposer que  $\nu$  est absolument continue par rapport à la mesure de Lebesgue  $d$ , de densité notée  $f$ :

$$\forall A \in \mathcal{B}(\mathbb{R}), \quad \nu(A) = \int_A f \, d$$

Par le théorème fondamental du calcul différentiel, cela est le cas si et seulement si  $F$  est absolument continue, si et seulement si  $F$  est presque partout dérivable, de dérivée  $F' = f$  avec en outre

$$F(x) = \int_{-\infty}^x f \, d \quad \text{p.p.}$$

- Nouvelle stratégie :

Il suffit alors d'estimer la densité  $f$ , qui est la dérivée de  $F$ , plutôt que  $F$  elle-même, cf. le théorème de Scheffé. On détaille cela :

### Le modèle de l'estimation de densité.

$X_1 \dots X_n$  iid selon une loi  $\nu \ll \lambda$ , de densité  $f$ :  $\frac{d\nu}{d\lambda} \stackrel{\text{not.}}{=} f$   
 $f$  est inconnue et on veut l'estimer; tout ce qu'on sait, c'est que  $f$  vit dans

$$\mathcal{P} = \left\{ g \in L^1(\lambda) : \int_{\mathbb{R}} g \, d\lambda = 1 \text{ et } g \geq 0 \right\}$$

(et  $\mathcal{P}$  est un très gros ensemble!).

Th: Soient  $\mu_1$  et  $\mu_2$  deux probabilités, avec  $\mu_1 \ll \lambda$  et  $\mu_2 \ll \lambda$ , on  
 (de Scheffé) l'on note  $f_1 = \frac{d\mu_1}{d\lambda}$  et  $f_2 = \frac{d\mu_2}{d\lambda}$ . Alors:

$$d_{\text{TV}}(\mu_1, \mu_2) = \int_{\mathbb{R}} (f_1 - f_2)_+ \, d\lambda = \frac{1}{2} \int_{\mathbb{R}} |f_1 - f_2| \, d\lambda$$

Preuve:  $B \stackrel{\text{not.}}{=} \{x \in \mathbb{R} : f_1(x) \geq f_2(x)\} \stackrel{\text{not.}}{=} \{f_1 \geq f_2\}$

Pour tout  $A \in \mathcal{B}(\mathbb{R})$ ,

$$\begin{aligned} \mu_1(A) - \mu_2(A) &= \int_A (f_1 - f_2) \, d\lambda \leq \int_{A \cap B} (f_1 - f_2) \, d\lambda \quad \text{par déf. de } B \\ &\leq \int_B (f_1 - f_2) \, d\lambda \quad \text{tp par déf. de } B \\ &= \int (f_1 - f_2)_+ \, d\lambda \end{aligned}$$

En passant à  $\sup_{A \in \mathcal{B}(\mathbb{R})}$ , on obtient la première égalité.

La seconde procède de ce que

$$0 = 1 - 1 = \int (f_1 - f_2) \, d\lambda = \int (f_1 - f_2)_+ \, d\lambda - \int (f_1 - f_2)_- \, d\lambda$$

$$\text{soit } \int (f_1 - f_2)_+ \, d\lambda = \int (f_1 - f_2)_- \, d\lambda;$$

or,  $\int |f_1 - f_2| d\mathbb{P} = \int (f_1 - f_2)_+ d\mathbb{P} + \int (f_1 - f_2)_- d\mathbb{P} = 2 \int (f_1 - f_2)_+ d\mathbb{P}$ ,  
ce qui conclut la preuve.

Stratégie: Estimer  $f$  de manière convergente pour la norme  $L^1$ :  
Construire  $\hat{f}_n \in \mathcal{F}$  à partir de  $X_1, \dots, X_n$  telle que

$$\|\hat{f}_n - f\|_1 \stackrel{\text{not.}}{=} \int |\hat{f}_n - f| d\mathbb{P} \xrightarrow{\mathbb{P}} 0$$

Remarque souvent utile: Comme  $\int |\hat{f}_n - f| d\mathbb{P}$  est toujours  $\leq 2$  quelles que  
soient  $\hat{f}_n$  et  $f$ , il vient:

$$\|\hat{f}_n - f\|_1 \xrightarrow{\mathbb{P}} 0 \iff \mathbb{E}[\|\hat{f}_n - f\|_1] \rightarrow 0.$$

(Ici, la convergence en probabilité et celle dans  $L^1(\mathbb{P})$  sont équivalentes pour notre critère.)

Estimation par fenêtres glissantes.

Son intuition repose sur le fait que  $f$  est la dérivée de  $F$ , et que  $F$  est bien approchée par la fonction de répartition empirique  $F_n$ :

$$\hat{f}_n = \frac{F_n(\cdot + h_n) - F_n(\cdot - h_n)}{2h_n} \quad \text{où } h_n > 0$$

est un paramètre à régler

Il faut:

- $h_n \searrow 0$  pour respecter l'intuition
- mais pas trop vite (on verra que  $nh_n \rightarrow +\infty$  est souhaitable) afin qu'il y ait suffisamment d'observations  $X_t$  dans chaque intervalle  $]x - h_n, x + h_n]$ :

$$\forall z \in \mathbb{R}, \quad \hat{f}_n(z) = \frac{1}{2nh_n} \sum_{t=1}^n \mathbb{1}_{\{X_t \in ]z - h_n, z + h_n]\}}$$

Lm:  $\hat{f}_n \in \mathcal{P}$

Preuve: Par Fubini-Tonelli,

$$\begin{aligned} \int_{\mathbb{R}} \hat{f}_n(x) dx &= \frac{1}{2nh_n} \sum_{t=1}^n \int \mathbb{1}_{\{X_t \in ]x-h_n, x+h_n]\}} dx \\ &= \frac{1}{2nh_n} \sum_{t=1}^n \underbrace{d([X_{t-h_n}, X_{t+h_n}[)}_{= 2h_n} \\ &= 1 \end{aligned}$$

Prop: Lorsque  $h_n \rightarrow 0$  et  $nh_n \rightarrow +\infty$ , on a en tout point  $x$  tel que  $F$  est dérivable en  $x$  (i.e., pour presque tout  $x$ ),

$$\mathbb{E}[(f(x) - \hat{f}_n(x))^2] \rightarrow 0.$$

Preuve: On commence par une décomposition biais-variance:

$$\begin{aligned} (*) \quad \mathbb{E}[(\hat{f}_n(x) - f(x))^2] &= \underbrace{\left( \mathbb{E}[\hat{f}_n(x)] - f(x) \right)^2}_{\text{cf. terme croisé nul}} + \mathbb{E}\left[ (\hat{f}_n(x) - \mathbb{E}[\hat{f}_n(x)])^2 \right] \\ &= \underbrace{\left( \mathbb{E}[\hat{f}_n(x)] - f(x) \right)^2}_{\text{biais}} + \underbrace{\text{Var } \hat{f}_n(x)}_{\text{variance}} \end{aligned}$$

Terme de biais:  $\mathbb{E}[\hat{f}_n(x)] = \frac{F(x+h_n) - F(x-h_n)}{2h_n} \rightarrow f(x)$  par hyp.  
 (cf. Hy,  $\mathbb{E}[F_n(y)] = F(y)$ ) ( $F'(x) = f(x)$  et  $h_n \rightarrow 0$ )

Terme de variance:  $2nh_n \hat{f}_n(x) = \sum_{t=1}^n \mathbb{1}_{\{X_t \in ]x-h_n, x+h_n]\}}$

suit une loi binomiale de paramètres  $n$  et:

$$p_{nx} = \mathbb{P}\{X_t \in ]x-h_n, x+h_n]\} = \frac{F(x+h_n) - F(x-h_n)}{2h_n}$$

$$\begin{aligned} \text{Donc } \text{Var } \hat{f}_n(x) &= \frac{1}{4n^2 h_n^2} \underbrace{\text{Var}(2nh_n \hat{f}_n(x))}_{= p_{nx}(1-p_{nx})} \leq \frac{1}{4nh_n} \frac{p_{nx}}{h_n} \sim \frac{1}{4nh_n} 2f(x) \\ &\rightarrow 0 \text{ car } nh_n \rightarrow +\infty \end{aligned}$$

Donc (\*) est la somme de deux termes tendant vers 0.

Th: Lorsque  $h_n \rightarrow 0$  et  $nh_n \rightarrow +\infty$ ,

$$\|\hat{f}_n - f\|_1 \xrightarrow{\mathbb{P}} 0$$

Preuve: Comme une remarque précédente l'indiquait, il suffit de montrer que

$$\begin{aligned} \mathbb{E}[\|\hat{f}_n - f\|_1] &\rightarrow 0 \\ &= 2 \mathbb{E}\left[\int (f - \hat{f}_n)_+ d\mu\right] \stackrel{\text{Fubini-Tonelli}}{=} 2 \int \varphi_n d\mu, \end{aligned}$$

où l'on a noté  $\varphi_n(x) = \mathbb{E}[(f(x) - \hat{f}_n(x))_+]$ .

Or,

- $\varphi_n \leq f \in L^1(\mu)$

- pour presque tout  $x$ ,

$$\varphi_n(x) \leq \sqrt{\mathbb{E}[(f(x) - \hat{f}_n(x))^2]} \xrightarrow{\text{Cauchy-Schwarz}} 0 \quad (\text{par la proposition})$$

D'où le résultat du théorème, par convergence dominée.

## Extensions et contenu du livre.

1 - Estimateurs à noyau : en notant  $K = \frac{1}{2} \mathbb{1}_{[-1,1]}$ , on a que l'estimateur par fenêtres glissantes correspond à

$$\hat{f}_{n,h_n} = \frac{1}{m \cdot h_n} \sum_{t=1}^n K\left(\frac{\cdot - X_t}{h_n}\right)$$

↳ D'autres choix de  $K$  sont possibles, qui conduisent à des estimateurs plus régulières (voir les illustrations graphiques)

↳ ... Mais il faudra prouver des théorèmes généraux de convergence!

OBJET DES PREMIÈRES SÉANCES

vs.  
OBJET DES SÉANCES SUIVANTES

2 - Choix de  $h_n$  : - beaucoup de suites  $(h_n)$  vérifient  $h_n \rightarrow 0$  et  $nh_n \rightarrow +\infty$   
↳ embarras du choix  
- en pratique,  $n$  est fixé (et non  $n \rightarrow +\infty$  !) et on choisit  $h_n$  en fonction de  $X_1, \dots, X_m$

qui est un paramètre d'utilisateur.

↳ Choix aléatoire  $h_n$  fondé sur les observations, qui conduit à un estimateur  $\hat{f}_{n,h_n}$ , à comparer aux autres choix de  $h$  possibles, qui auraient conduit aux estimateurs  $\hat{f}_{n,h}$  ; on veut faire presque aussi bien que le meilleur de ces derniers.

On va établir l'inégalité exacte suivante :

pour un choix  $\hat{h}_n$  précis de  $h_n$  en fonction de  $X_1, \dots, X_m$  uniquement :

$$\frac{\sup_{f \in \mathcal{F}} E[\|\hat{f}_{n,\hat{h}_n} - f\|_1]}{\inf_{h>0} E[\|\hat{f}_{n,h} - f\|_1]} \leq 3 + o(1)$$

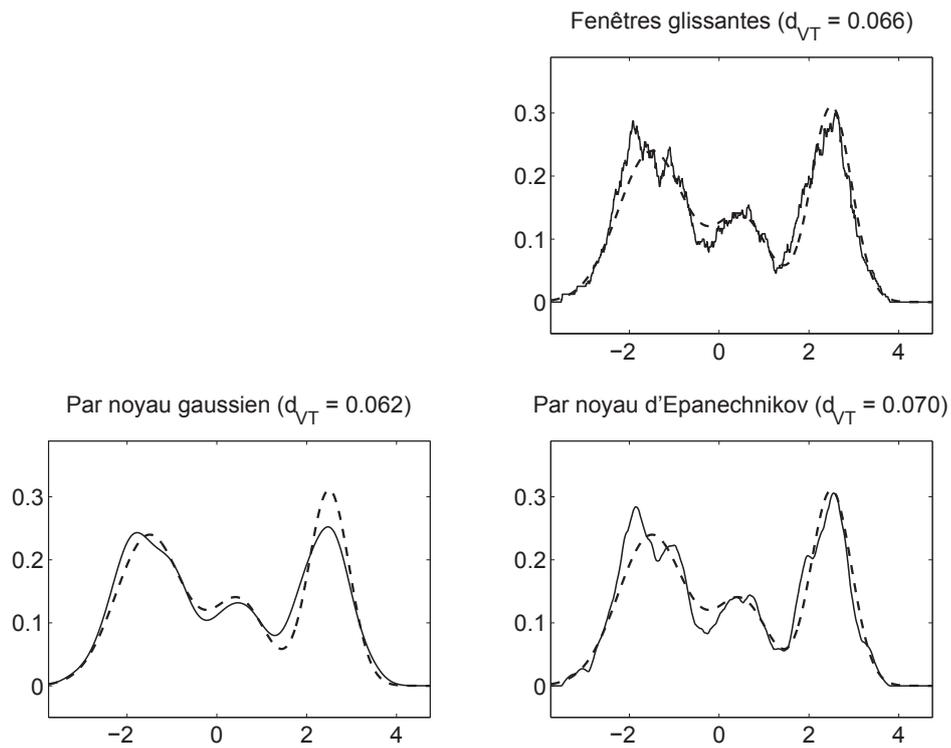


Figure 1. Une réalisation des estimations pour le mélange de lois gaussiennes,  $n = 400$  et  $h = 0.30$ . On précise entre parenthèses les distances en variation totale réalisées.

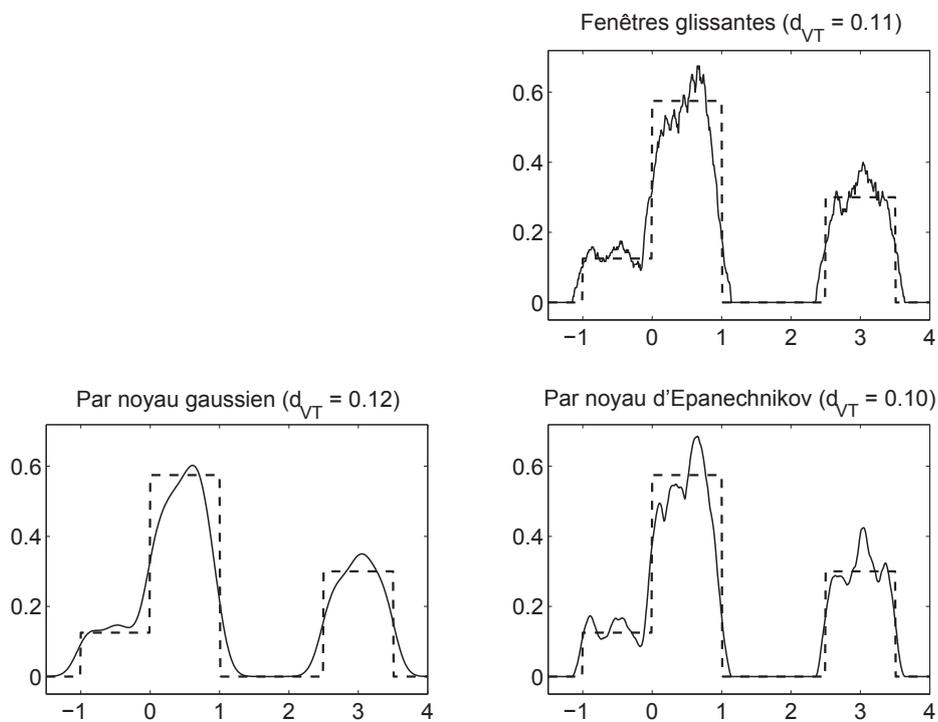


Figure 2. Une réalisation des estimations pour le mélange de lois uniformes,  $n = 400$  et  $h = 0.15$ . On précise entre parenthèses les distances en variation totale réalisées.

### 3. Ingrédients mis en jeu :

- des techniques d'analyse élémentaires (convolution dans  $L^1(\mathbb{R}^d)$ ) pour la convergence des estimateurs par noyau
- de la concentration de sommes de variables aléatoires
- un peu de combinatoire (dimension de Vapnik, lemme de Sauer)

PLAN : DL00 sera notre référence principale

Deux références secondaires détaillant certains points délicats sont :

- Dev85 (Devroye : « A course in density estimation »)
- RS09 (Riviere et Stoltz : « Statistique en Action »)

*(ces trois ouvrages sont disponibles à la bibliothèque)*

Séance 1 : Convergence des estimateurs à noyau à suite  $(h_n)$  fixé

(7 octobre)

(techniques d'analyse)

DL00, § 9.1 - 9.3

RS09, § 24.6

Séance 2 : Inégalités de concentration pour (essentiellement) des

(14 octobre)

sommes de variables aléatoires iid

DL00, § 2.1 - § 2.3

RS09, § 26.2.2

Séance 3 : Convergence des estimateurs à noyau : choix  $(H_n)$

(21 octobre)

à partir des données  $X_1, \dots, X_n$  ; convergence ps

(et pas seulement  $L^1$  ou en probabilité)

DL00, Th 9.3 ;  
plus: § 2.4 et § 9.4

Dev85, § 3.6

Séance 4 : Un mot sur les vitesses de convergence : bornes supérieures uniquement

(28 octobre)

(encore une séance avec pas mal d'analyse)

DL00 § 9.8

Il s'agit de faire simplement une mise en perspective ; voir aussi RS09, § 32.1.4  $\left\{ \begin{array}{l} \text{+ éléments au choix dans} \\ \text{le § 9.7 et les exercices (ex 9.24)} \end{array} \right.$

Séance 5: Méthodes et techniques combinatoires DL00 § 3.1 et § 3.3  
(4 novembre)

Séance 6: Comment choisir entre un nombre fini de densités? Vers une première inégalité oracle, à améliorer. DL00, Théorèmes: 6.3, 6.4 et 10.1

Séance 7: Illustrations de l'inégalité oracle sur la classe des estimateurs de forme additive DL00, § 10.2

Séance 8: L'inégalité oracle finale. DL00 § 11.1 et § 11.2 et des éléments du § 11.3 si besoin.

En cas d'affluence, deux exposés supplémentaires sont disponibles:

Séance 5 bis: Méthodes et techniques combinatoires, suite DL00 § 4.1, § 4.2 et § 4.4

Séance 6 bis: Autres illustrations de la première inégalité oracle Chap 8 (choisir quelques exemples, ne pas tous les traiter)