

TP3 - Statistique Non Paramétrique

Estimation non paramétrique d'une régression

Le rapport est à rendre avant le mercredi prochain (20 Avril).

Exercice 1 L'estimateur régressogramme

- (a) Construire l'estimateur régressogramme (à classes de même taille) d'une fonction de régression m définie sur $[0, 1]$ en complétant le code ci-dessous :

```
> # construire l'estimateur régressogramme comme une fonction de N, de l'échantillon (X_i, Y_i)
> # ou N est le nombre de classes / intervalles / bins
> # et x est un point dans [0,1]
> estimateur_reg <- fonction(x, N, echX, echY){
+   #assertation
+   stopifnot(length(echX)==length(echY))
+   n <- length(echX)
+   # découper [0, 1] en N classes de même taille
+   breakpoints <- seq(0, 1, length.out = N + 1)
+   # calculer la valeur de l'estimateur régressogramme en point x: <<m_chap_x>>
+   for (i in 1:N){
+     # réperer la bonne classe <<j>> ou se trouve x
+     if(x >= breakpoints[i] & x < breakpoints[i+1]){
+       # calculer la moyenne des Y_i dont X_i tombent dans la classe <<j>>
+
+
+       is_echX_in_classj <- echX >= breakpoints[i] & echX < breakpoints[i+1]
+
+       echY_in_classj <- echY[is_echX_in_classj]
+
+       m_chap_x <- mean(echY_in_classj)
+
+
+     }
+   }
+ }
```

```
+ return(m_chap_x)
+ }
```

(b) **Simulation : l'évaluation d'une méthode sur les données idéales.**

Utiliser le code ci-dessous pour

- simuler un $n = 1000$ - échantillon $(X_1, Y_1), \dots, (X_n, Y_n)$ i.i.d. de (X, Y) dont la vraie relation est donné par $Y = m(X) + \varepsilon$. On prends $X_i \stackrel{i.i.d.}{\sim} \mathcal{U}(0, 1)$, $m(x) = \sin(2\pi x^2)^2$, et $(\varepsilon_1, \dots, \varepsilon_n)$ i.i.d de $\varepsilon \sim \mathcal{N}(0, 0.5^2)$.

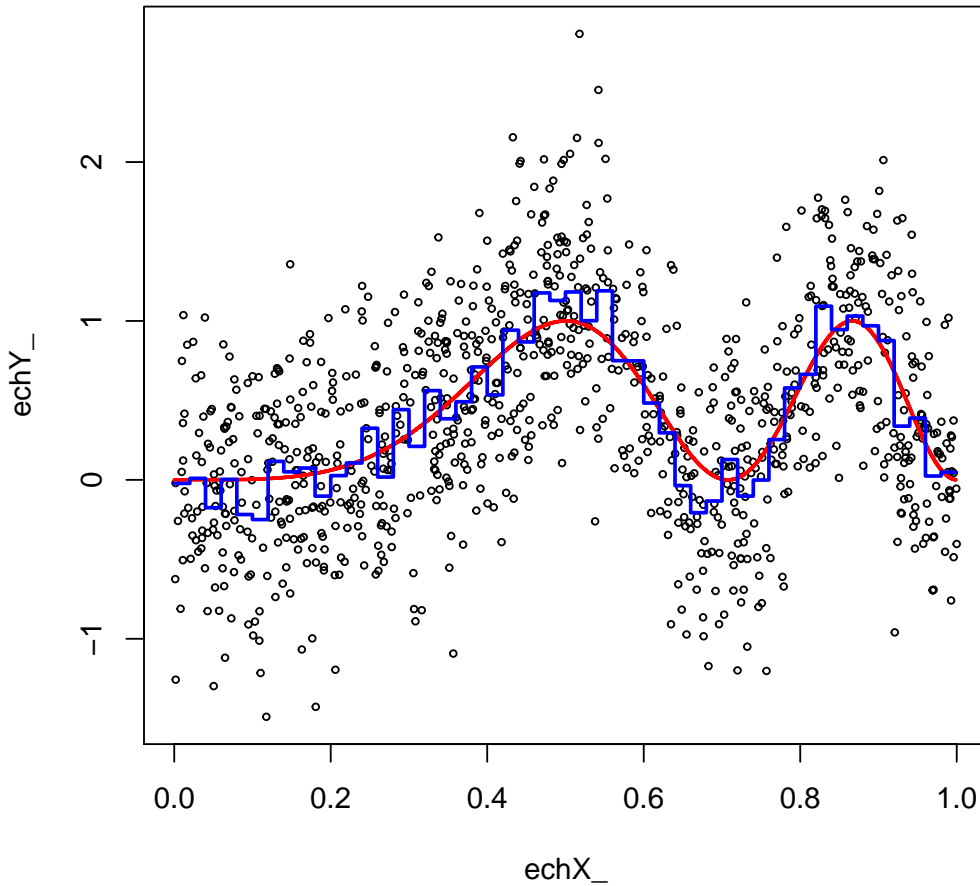
Le but de la simulation est de vérifier si l'estimateur proposé peut trouver la vraie relation m depuis les données idéales (les hypothèses sont vérifiées). Sioui, l'estimateur est alors valable pour les données réelles.

- appliquer l'estimateur sur les données simulées avec $N = 10, 30, 50$,

- tracer la fonction estimée sur $[0, 1]$ (superposer les données simulées et la vraie relation m à estimer).

Quel N est le meilleur ?

```
> # le modèle de regression est vrai pour les données simulées
> set.seed(1)
> echX_ <- runif(1000)
> echEpsilon <- rnorm(1000, 0, 0.5)
> echY_ <- sin(2*pi*echX_^2)^2 + echEpsilon
> # appliquer et tracer l'estimateur sur les données simulées
> x_list <- seq(0,0.999,length.out = 1000)
> estimation_reg <- sapply(x_list, function(x_){
+                                     return(estimateur_reg(x_,50,echX_,echY_
+                                     )
+                                     )
+                                     }
+                                     )
> plot(echX_, echY_, cex = 0.5)
> lines(x_list,sin(2*pi*x_list^2)^2, type="s", col = "red", lwd=2)
> lines(x_list,estimation_reg, type="s", col = "blue", lwd=2)
```



$N = 30$ est le meilleur.

Exercice 2 L'estimateur à noyau

- (a) Construire un estimateur à noyau (la densité de X est inconnue) avec le noyau Gaussien en complétant le code ci-dessous.

Rapel de la définition :

Si $\mathbf{K} : \mathbb{R} \rightarrow \mathbb{R}$ est un noyau statistique d'ordre 1, *l'estimateur de Nadaraya-Watson* vaut

$$\hat{m}_{NW}(x) = \begin{cases} \frac{\sum_{i=1}^n \mathbf{K}\left(\frac{X_i - x}{h}\right) Y_i}{\sum_{k=1}^n \mathbf{K}\left(\frac{X_k - x}{h}\right)} & \text{si } \sum_{k=1}^n \mathbf{K}\left(\frac{X_k - x}{h}\right) \neq 0 \\ 0 & \text{sinon.} \end{cases} \quad (0.1)$$

```
> #construire l'estimateur à noyau Gaussien comme une fonction de h, de l'échantillon (X, Y)
> estimateur_noyau_Gaussien <- fonction(x, h, echX, echY){
+   #assertation
+   stopifnot(length(echX)==length(echY))
```

```

+ # calculer la valeur de l'estimateur à noyau en point x: <<m_chap_x>>
+
+ int_K <- (echX - x)/h
+ m_chap_x <- sum(dnorm(int_K)*echY)/sum(dnorm(int_K))
+
+
+ return(m_chap_x)
+ }

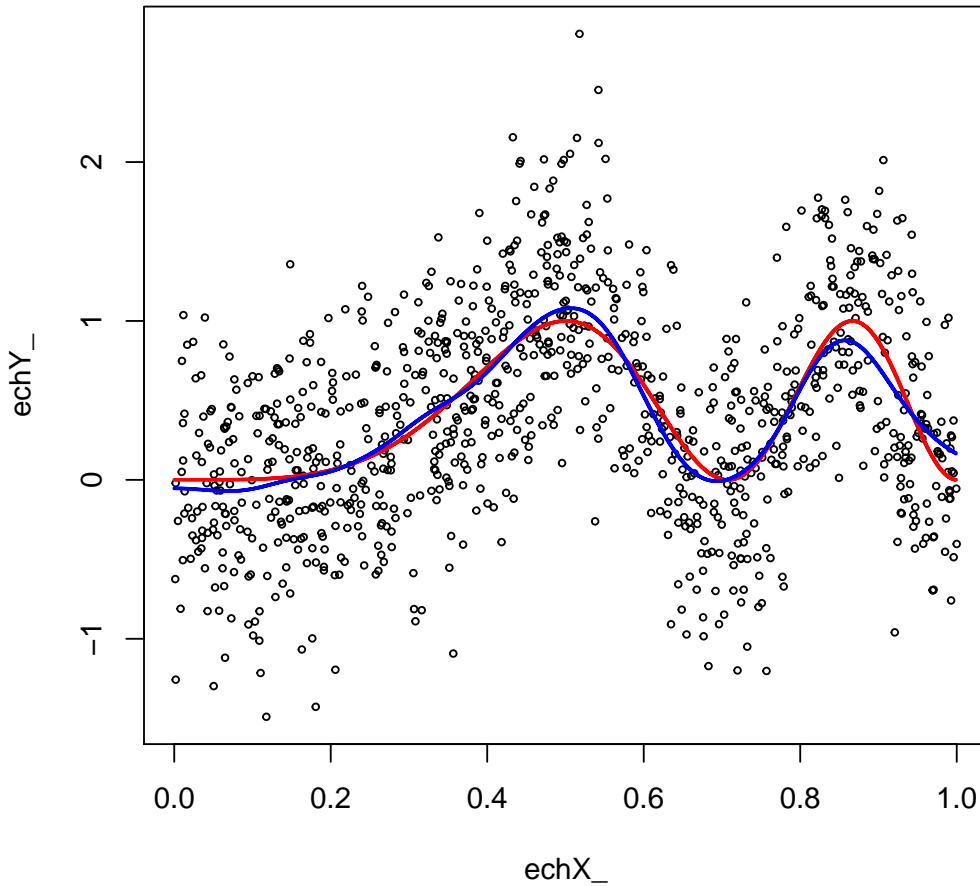
```

- (b) Appliquer l'estimateur à noyau Gaussien à l'échantillon de l'exercice 1 avec 3 valeurs de h à votre choix qui donnent respectivement sous-lissage, sur-lissage, et ni l'un ni l'autre dans leurs estimations. On notera cette dernière h^* dans la suite.

```

> x_list <- seq(0,0.999,length.out = 1000)
> estimation_noyau_Gaussien <-
+ sapply(x_list, function(x){return(estimateur_noyau_Gaussien(x,h=0.04,echX_, echY_))})
> plot(echX_, echY_, cex = 0.5)
> lines(x_list,sin(2*pi*x_list^2)^2, type="s", col = "red", lwd=2)
> lines(x_list,estimation_noyau_Gaussien, type="s",col="blue", lwd=2)

```



0.01 : sous-

lissage, 0.1 : sur-lissage

Est-ce que vous avez pu trouver le parfait h , il semble qu'il y a toujours une partie de estimation qui colle pas sur la vraie courbe ? e.g. $h = 0.05$ sauf que le bout à droite, l'estimation est déjà parfaite, ici l'estimation sur le bout à droite est moins variable que la vraie, ctd on a sur-lissage ici (h est trop grand), mais si on baisse h à e.g. 0.04, on a sous-lissage tout de suite sur les autres parties de l'intervalle.

D'où vient le probleme ? quelle est la difference niveau données entre le bout et le centre ?

L'estimateur manque les données quand on s'approche des bouts de l'intervalle. Donc l'estimation sur 2 bouts de l'intervalle demande une plus petite valeur h , mais puisqu'on utilise qu'une seule valeur donc le parfait h n'existe pas, on peut par exemple déterminer h en regardant plutôt l'estimation vers le centre si l'intervalle est assez long.

- (c) Retrouver votre estimation avec celui donné par R en utilisant les fonctions `ksmooth` et `locpoly` pour $h = h^*$, à l'aide du code ci-dessous.

```
> h_star <- 0.04
> # le package utilise une autre déf du bandwidth, en fonction de quartiles de la Gaussian
> # qui est égale à 2.67*h (avec h de notre définition)
```

```

> # plus d'explication:
> # https://stats.stackexchange.com/questions/396324/what-does-bandwidth-in-kernel-regres
> estimation_ksmooth <- ksmooth(echX_, echY_, kernel = "normal", bandwidth = 2.67*h_star,
> plot(echX_, echY_, cex = 0.5)
> lines(x_list, estimation_noyau_Gaussien,col="green", lwd=3)
> lines(estimation_ksmooth,col="blue",type='s', lwd=1)

```

La méthode d'estimation de Nadaraya-Watson est un cas particulier de la technique des polynômes locaux qui est réalisée dans R avec la fonction `locpoly`. L'estimateur de Nadaraya-Watson=l'estimateur est lié au degré zéro et s'obtient ici :

```

> library(KernSmooth)
> estimation_locpoly <- locpoly(echX_, echY_,degree=0,bandwidth=h_star,gridsize=1000,range)
> plot(echX_, echY_, cex = 0.5)
> lines(x_list, estimation_noyau_Gaussien,col="green", lwd=3)
> lines(estimation_locpoly,col="blue",type='s', lwd=1)

```

- (d) **Observer l'effet local du lissage.** Ajouter $(0.5, -5)$ dans l'échantillon, et puis appliquer l'estimateur à noyau sur l'échantillon augmenté avec $h = h^*$. Quelle est la différence entre cette dernière estimation et la précédente?

```

> estimation_noyau_Gaussien_aug <-
+   sapply(x_list, fonction(x){
+     return(estimateur_noyau_Gaussien(
+       x_,h=h_star, c(echX_,0.5), c(echY_, -5)
+     )))
> plot(echX_, echY_, cex = 0.5)
> lines(x_list,estimation_noyau_Gaussien_aug, type="s",col="red", lwd=2)
> lines(x_list,estimation_noyau_Gaussien, type="s",col="blue", lwd=2)

```

Dans un voisinage de la nouvelle observation $(0.5, -5)$, les valeurs d'estimation sont tous biaisées. De plus, plus loins de $x = 0.5$, moins de changement sur les valeurs d'estimation. Parce que le poids pour cette nouvelle observation $(0.5, -5)$ diminue quand on s'en éloigne.

Exercice 3 Sélection de la fenêtre h : la validation croisée (Leave One Out)

On va maintenant chercher une approche qui permet de choisir l'hyperparamètre h en autonomie (dans la pratique), et puis la programmer.

On d'abords cherche une métrique qui évalue la performance de l'estimateur à noyau \hat{m} appliqué sur un jeu de données $(X_1, Y_1), \dots, (X_n, Y_n)$. Un choix courant en pratique est la somme des carrés des résidus (SCR) :

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{m}_h(X_i))^2.$$

Donc, la valeur de h qui minimise la SCR est optimale. Mais en effet le minimiseur est toujours $h = 0$ pour la formule ci-dessus, car \hat{m}_h est calculé avec Y_i déjà (plus d'explication, cf. Section 4.3.2 <https://bookdown.org/egarpor/NP-UC3M/kre-i-bwd.html#kre-i-bwd-cv>), donc on considère au

lieu son variant

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{m}_{h,-i}(X_i))^2, \quad (0.2)$$

où $\hat{m}_{h,-i}(x)$ est l'estimateur de Nadaraya-Watson sur les n couples privés du couple (X_i, Y_i) .

Comme la validation croisée utilisée dans le cadre de l'estimation de densité, on cherche à lier $\hat{m}_{h,-i}(x)$ à $\hat{m}_h(x)$ pour faciliter le code.

Pour ceci, on note $\omega_j(x)$ le poids $\frac{\mathbf{K}\left(\frac{X_j - x}{h}\right)}{\sum_{k=1}^n \mathbf{K}\left(\frac{X_k - x}{h}\right)}$, donc $\hat{m}_h(x) = \sum_{j=1}^n \omega_j(x) Y_j$.

D'autre part, on note $\omega_j^{-i}(x)$ le poids $\frac{\mathbf{K}\left(\frac{X_j - x}{h}\right)}{\sum_{k=1, k \neq i}^n \mathbf{K}\left(\frac{X_k - x}{h}\right)}$, donc $\hat{m}_{h,-i}(x) = \sum_{j=1, j \neq i}^n \omega_j^{-i}(x) Y_j$.

D'après le calcul direct, on peut trouver

$$\omega_j^{-i}(x) = \frac{\omega_j(x)}{1 - \omega_i(x)}. \quad (0.3)$$

Injecter Equation 0.3 dans $\hat{m}_{h,-i}(x)$, on a trouvé ainsi le lien désiré

$$\hat{m}_{h,-i}(x) = \frac{\hat{m}_h(x) - \omega_i(x) Y_i}{1 - \omega_i(x)}$$

La SCR (0.2) est devenue

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \left(\frac{Y_i - \hat{m}_h(X_i)}{1 - \omega_i(X_i)} \right)^2.$$

- (a) Créer une fonction de R qui prend la fenêtre h et l'échantillon $(X_1, Y_1), \dots, (X_n, Y_n)$ comme les entrées et sort $\hat{R}(h)$, en complétant le code ci-dessous.

```
> R_chap <- fonction(h, echX, echY){
+   #calculer le vecteur numérateur: (Y1-m_chap_h(X1), Y2-m_chap_h(X2), ..., Yn-m_chap_h(Xn))
+
+   vect_numérateur <- echY - sapply(echX,
+     fonction(x){
+       return(estimateur_noyau_Gaussien(x, h, echX, echY))
+     }
+   )
+
+   #calculer le vecteur dénominateur: (1-w1(X1), 1-w2(X2), ..., 1-wn(Xn))
+   wi_Xi <- fonction(Xi){
+     int_K <- (echX - Xi)/h
+     val <- dnorm(0)/sum(dnorm(int_K))
+     return(val)
+   }
+ }
```

```

+ vect_dénominateur <- 1-sapply(echX, wi_Xi)
+
+ return( mean( vect_numérateur^2 / vect_dénominateur^2 ) )
+ }

```

- (b) Tracer la courbe de $\hat{R}(h)$ en fonction de h , avec l'échantillon donné dans les exercices précédents, puis repérer son minimum.

```

> h_list <- seq(0.01, 0.1, 0.001)
> val_R_chap <- sapply(h_list, function(h_) R_chap(h_, echX_, echY_))
> plot(h_list, val_R_chap, type = "s")
> # h optimal
> h_list[which.min(val_R_chap)]
[1] 0.024

> estimation_noyau_Gaussien <-
+ sapply(x_list, function(x_){return(estimateur_noyau_Gaussien(x_,h=0.024,echX_, echY_))})
> plot(echX_, echY_, cex = 0.5)
> lines(x_list,sin(2*pi*x_list^2)^2, type="s", col = "red", lwd=2)
> lines(x_list,estimation_noyau_Gaussien, type="s",col="blue", lwd=2)

```