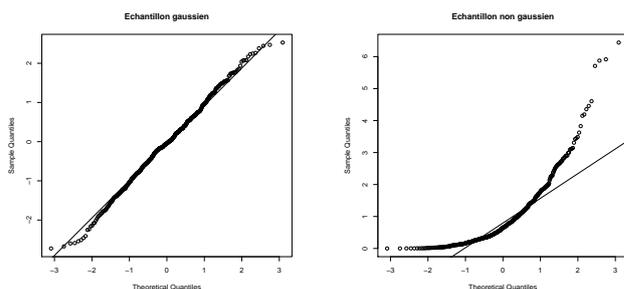


TP 2 : Régression non-paramétrique

Préambule. Dans ce TP, nous allons implémenter des procédures de régression non-paramétriques par noyaux et par polynômes locaux. Vous aurez besoin de charger (fonction `library`) les packages `locfit` et `np`.

On rappelle les fonctions suivantes :

- `par(mfrow=c(k,1))` permet d'afficher $k \times l$ graphiques sur une même figure.
- `lines()` (resp. `points()`) permet d'ajouter une courbe (resp. des points) sur un graphe existant.
- `qqnorm(X)` permet de comparer graphiquement la distribution d'un échantillon X avec une distribution normale (si les points sont approximativement alignés, X est gaussien). `qqline(X)` permet d'ajouter la droite passant par le premier et troisième quartile :
`qqnorm(X)`
`qqline(X)`



`qqplot(X,Y)` permet de comparer graphiquement les distributions de deux échantillons X et Y . Si les distributions sont identiques, les points sont approximativement alignés.

- `seq(a,b,l=m)` fournit un vecteur de m points équidistants entre a et b .

Données

Nous allons travailler sur 3 jeux de données :

- **Data A.** Les données `cps71` du package `np` fournissent le logarithme du salaire et l'âge pour 205 canadiens.
- **Data B.** Les données du data.frame `Alcool` fournissent la consommation d'alcool en g/jour pour 312 individus (variable `Alcohol.consumption`) et leur durée de vie divisée par la durée de vie moyenne (variable `Life`).
- **Data C.** Les données du data frame `Hormone` fournissent la dose d'hormone administrée au patient (variable `Hormone`) et le taux d'endocrine résultant dans le sang (variable `Endocrine`) lors d'une étude clinique.

Rq : Les données Data B et C ont été simulées, mais les distributions ont été choisies de manière réaliste sur la base de travaux scientifiques.

(I) Méthode de Nadaraya-Watson et Polynômes locaux

Dans cette section, vous allez vous familiariser avec des estimateurs de régression NP, en vous appuyant sur le jeu de données Data A.

I-1-a) Charger le jeu de données `cps71` (commande `data(cps71)`). Afficher le début de la matrice de données `cps71` à l'aide de la fonction `head`.

I-1-b) On veut étudier la fonction de régression de la variable X_i égale à l'âge (colonne `age`) sur la variable Y_i égale au logarithme du salaire (colonne `logwage`). Définir les vecteurs `X` et `Y` correspondant de longueur $n = 205$, et tracer les points d'observations $(X_i, Y_i)_{i=1, \dots, n}$.

I-1-c) Définir une grille `x0` de 1000 points régulièrement espacés sur l'intervalle $[\min(X_1, \dots, X_n), \max(X_1, \dots, X_n)]$.

I-2) **Régression par noyaux (estimateur de Nadarya Watson)**. L'estimateur de Nadaraya-Watson est implémenté dans la fonction `npreg` du package `np`. L'argument principal est la formule `Y~X` où `Y` est la variable réponse et `X` la variable explicative.

```
mod.ker <- npreg(Y~X)
```

La commande `names(mod.ker)` permet d'afficher le nom de l'ensemble des variables fournies par `npreg`.

La valeur de l'estimateur $\hat{r}(x)$ pour x dans un vecteur `X0` choisi par l'utilisateur est obtenue en spécifiant l'argument `exdat` dans la fonction `npreg`. Les valeurs $(\hat{r}(x), x \in X_0)$ sont alors contenues dans la variable `mean`.

```
X0 <- seq(30, 50, l=100)
mod.kern <- npreg(Y~X, exdat=X0)
rhatX0 <- mod.kern$mean
plot(X0, rhatX0)
```

Si on ne spécifie pas la valeur de `exdat`, `npreg` calcule l'estimateur aux points du vecteur `X`.

La fenêtre optimale est calculée par une validation croisée (interne) mais peut également être choisie par l'utilisateur. Dans ce TP, nous travaillerons avec la valeur déterminée par `npreg`.

Question. Sur un même dessin, tracer les points d'observations $(X_i, Y_i)_{i=1, \dots, n}$ ainsi que l'estimateur \hat{r} par noyau avec la fenêtre par défaut, calculé aux points de `x0`. On notera `rkern` le vecteur contenant les valeurs $(\hat{r}(x), x \in \mathbf{x0})$.

I-3) **Estimateurs par polynômes locaux**. La fonction `locfit` du package `locfit` réalise la régression par polynôme locaux. L'argument principal est la formule `Y~lp(X)` où `Y` est la variable réponse et `X` la variable explicative.

```
mod.lp <- locfit(Y~lp(X))
```

Cet estimateur comporte deux paramètres de régularisation (le degré maximum du polynôme, et une quantité qui contrôle la fonction de poids). Ces paramètres comportent des valeurs par défaut, et peuvent également être contrôlés par l'utilisateur. Dans ce TP, nous travaillerons avec les valeurs par défaut.

La valeur de l'estimateur $\hat{r}(x)$ pour x dans un vecteur `X0` choisi par l'utilisateur est obtenue dans un second temps à l'aide de la fonction `predict` qui requiert (au moins) deux arguments :

un "modèle" de régression, et un ensemble de valeur où la fonction de regression doit être appliquée :

```
mod.lp <- locfit(Y~lp(X))
X0 <- seq(30,50,l=100)
rhatX0 <- predict(mod.lp, X0)
```

Question. Sur un même dessin, tracer les points d'observations $(X_i, Y_i)_{i=1, \dots, n}$ ainsi que l'estimateur \hat{r} par polynômes locaux, calculé aux points de $\mathbf{x0}$. On notera $\mathbf{r1p}$ le vecteur contenant les valeurs $(\hat{r}(x), x \in \mathbf{x0})$.

I-4) La fonction `lowess` implémente un estimateur de régression qui mélange les k plus proches voisins et la régression par polynômes locaux de degré 1. Cette fonction est très souvent utilisée pour un examen visuel. La régularisation est contrôlée par le paramètre f (plus f est grand, plus l'estimateur est lissée); par défaut $f=2/3$.

```
plot(X,Y)
lines(lowess(X,Y), col='red')
lines( lowess(X,Y, f=0.1), col='blue')
```

(II) Validation croisée sur le jeu de données Data B

II-1-a) Charger le jeu de données Data B. On veut étudier la fonction de régression de la variable X_i égale à la consommation d'alcool sur la variable Y_i égale à l'âge de décès normalisé. Définir les vecteurs \mathbf{X} et \mathbf{Y} correspondant de longueur $n = 312$, et tracer les points d'observations $(X_i, Y_i)_{i=1, \dots, n}$.

II-1-b) Définir une grille $\mathbf{x0}$ de 1000 points régulièrement espacés sur l'intervalle $[\min(X_1, \dots, X_n), \max(X_1, \dots, X_n)]$

II-2-a) Calculer le vecteur \mathbf{rkern} contenant les valeurs de l'estimateur par noyau aux points de $\mathbf{x0}$, et le vecteur $\mathbf{r1p}$ contenant les valeurs de l'estimateur par polynômes locaux aux points de $\mathbf{x0}$.

II-2-b) Tracer sur un même graphique :

- Les points $(X_i, Y_i)_{i=1, \dots, n}$
- L'estimateur par noyau, en rouge (argument `col='red'`)
- L'estimateur par polynômes locaux, en bleu (argument `col='blue'`)

II- 3) Dans cette question, nous allons voir pas-à-pas comment calculer l'erreur de validation croisée 10-fold, et tracer les valeurs $(\hat{Y}_i)_{i=1, \dots, n}$ estimées par VC en fonction des valeurs observées $(Y_i)_{i=1, \dots, n}$.

II-3-a) **Découpage de l'échantillon** : $\{1, \dots, n\}$: $n = 312$ n'est pas divisible par 10, nous allons donc diviser $\{1, \dots, n\}$ en 10 sous-ensembles disjoints de taille 31 ou 32. Soit

```
J1 <- sample(rep(1:10, l=n))
```

La commande `table(J1)` donne les effectifs de chaque valeur dans le vecteur $\mathbf{J1}$. Le vecteur $\mathbf{J1}$ contient 32 fois les valeurs 1 et 2, et 31 fois les valeurs 3,4,...,10, dans un ordre aléatoire. Pour tout $\ell = 1, \dots, 10$, I_ℓ est défini par les positions des valeurs ℓ dans le vecteur $\mathbf{J1}$.

Exple :

```
I2 <- which(J1==2)
```

II-3-b) Soit $l=1$. L'ensemble d'apprentissage l correspond aux observations telles que $J1 \neq 1$, et l'ensemble de validation aux observations telles que $J1 = 1$

- Calculer l'estimateur de régression par polynômes locaux à partir de l'ensemble d'apprentissage l :
 - variable explicative : `X[which(J1!=1)]`
 - variable réponse : `Y[which(J1!=1)]`
- Calculer le vecteur `Yhat` des valeurs prédites par VC pour l'ensemble de validation l (variables explicatives `X[which(J1==1)]`)
- Tracer les valeurs prédites `Yhat` en fonctions des valeurs observées `Y[which(J1==1)]` pour les observations de l'ensemble de validation l .

II- 3-c) Nous allons appliquer la même procédure pour tous les $l = 1, \dots, 10$:

- Définir un vecteur `Yhat.lp` de longueur n . Ce vecteur servira à stocker les valeurs prédites pour $l = 1, \dots, 10$.
- Pour tout $l \in \{1, \dots, 10\}$: calculer le vecteur `Yhat` des valeurs prédites par VC pour les observations de l'ensemble de validation l $\{\hat{Y}_i, i \in I_l\}$, et les stocker dans le vecteur `Yhat.lp` à l'aide de la commande :
`Yhat.lp[which(J1==1)] <- Yhat`

II-3-d) Tracer les valeurs prédites par VC en fonction des valeurs observées $(Y_i, \hat{Y}_i)_{i=1, \dots, n}$, ainsi que la droite $y = x$. Calculer la corrélation de Pearson entre $(Y_i)_{i=1, \dots, n}$ et $(\hat{Y}_i)_{i=1, \dots, n}$. Que pensez-vous de la qualité de prédiction ? Basé sur ces données, pensez-vous qu'il y a un lien entre consommation d'alcool et mortalité ?

II-3-e) Calculer l'erreur de VC `Ecv.lp`

II-4) La véritable fonction de régression (utilisée pour simuler les données) est :

$$r(x) = 10^{-4}(x + 40)^3 \exp\left(-\frac{x + 40}{20}\right)$$

Sur un même graphe, tracer les observations $(X_i, Y_i)_{i=1, \dots, n}$, la vraie fonction r et l'estimateur par polynômes locaux. Que pensez vous de la qualité d'estimation ?

II-5) Comment expliquez-vous qu'on puisse avoir une bonne qualité d'estimation mais une mauvaise prédiction ?

(III) Validation croisée sur le jeu de données Data C

Dans cette section, nous allons analyser le jeu de données Data C. Les questions suivantes sont très similaires à celles de la section précédente. Il peut-être judicieux de réutiliser les scripts...

III-1) Charger le jeu de données Data C. On veut étudier la fonction de régression de la variable X_i égale à la dose d'hormone prescrite sur la variable Y_i égale au taux d'endocrine

mesuré. Définir les vecteurs \mathbf{X} et \mathbf{Y} correspondant de longueur $n = 100$, et tracer les points d'observations $(X_i, Y_i)_{i=1, \dots, n}$.

Définir une grille $\mathbf{x0}$ de 1000 points régulièrement espacés sur l'intervalle $[\min(X_1, \dots, X_n), \max(X_1, \dots, X_n)]$

III-2) Tracer sur un même graphique :

- Les points $(X_i, Y_i)_{i=1, \dots, n}$
- L'estimateur par noyau, en rouge (argument `col='red'`)
- L'estimateur par polynômes locaux, en bleu (argument `col='blue'`)

III-3) Calculer les valeurs prédites par VC (10-fold) pour l'estimateur par polynômes locaux, et tracer les valeurs prédites en fonctions des valeurs observées. Que pensez-vous de la qualité de prédiction ?

III-4) Un simple coup d'oeil aux données montre qu'un modèle linéaire n'est pas adapté. Néanmoins, à titre d'exercice, on va comparer les résultats ci-dessus avec les prédictions obtenues par un modèle linéaire.

Prédiction dans un modèle linéaire aux points d'un vecteur $\mathbf{X0}$, pour un modèle estimé à partir d'un vecteur de variables explicatives \mathbf{X} et d'un vecteur de variables réponse \mathbf{Y} :

```
mod <- lm(Y~X)
co <- coef(mod)
rhatX0 <- co[1] + co[2]*X0
```

Calculer les valeurs `Yhat.lm` prédites par VC pour un modèle linéaire. Tracer ces valeurs prédites en fonction des valeurs observées et conclure sur les capacités prédictives du modèle. Comparer les erreurs de VC avec l'estimateur par polynômes locaux et l'estimateur linéaire.