

Exercice 1 (K-Nearest Neighbours):

Consider the KNN rule with equal weights for classification. One observes n i.i.d. replications $D_n = \{(X_i, Y_i), i \leq n\}$ of a random pair (X, Y) with $X \in \mathbb{R}^d$ and $Y \in \{0, 1\}$. Recall that the KNN classifier given the dataset D_n is

$$g_{k,n}(x) = \begin{cases} +1 & \text{if } \sum_{i=1}^k Y_{(i)}(x) \geq 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

Where $Y_{(i)}$ is the label of the i^{th} nearest neighbour $X_{(i)}$ of the point x . Recall also that the regression function is defined as $\eta(x) = \mathbb{P}[Y = 1 | X = x]$. We assume here that X is uniformly distributed over $[0, 1]$ and that $\eta(x) = 1/3$ for all $x \in [0, 1]$.

1. Compute the joint probability $\mathbb{P}[X \in [a, b], Y = 1]$ for all $0 < a < b$.
2. Give the expression for the Bayes classifier g^* For a generic regression function $\eta(x)$ and for this particular problem. Give its error risk $\mathbb{P}[g^*(X) \neq Y]$ (the Bayes error) in general and in this particular case.
3. Consider the training set

$$(X_1 = 0.4, Y_1 = 0), (X_2 = 0.2, Y_2 = 0), (X_3 = 0.7, Y_3 = 1)$$

Compute the k-NN classifier (for all x) in the case $k = 1$ and $k = 3$.

4. Compute in each case the expected classification error $R(g_{knn}) = \mathbb{P}[g_{knn}(X) \neq Y]$ for a new observation (X, Y) . Compare with the Bayes error
5. Consider now another model : X is again uniformly distributed over $[0, 1]$ but now $\eta(x) = \mathbb{1}_{(1/2, 1]}(x)$. Repeat questions 1,2,3,4 with the same dataset as in Q3.

$$\begin{aligned} 1) \mathbb{P}(X \in [a, b], Y = 1) &= \mathbb{E}[\mathbb{P}(Y = 1, X \in [a, b] | X)] \\ &= \mathbb{E}[\mathbb{P}(Y = 1 | X) \mathbb{1}_{[a, b]}(x)] \\ &= \int_a^b \frac{1}{3} dx = \frac{b-a}{3}. \end{aligned}$$

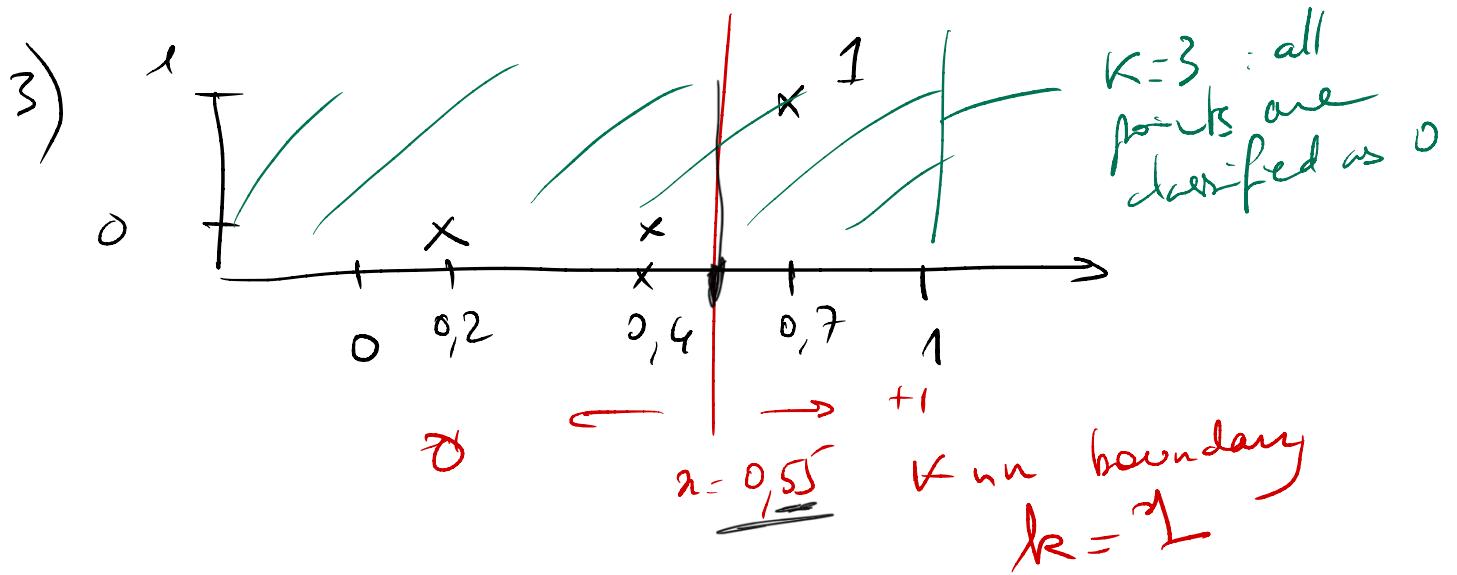
$$2) \text{ if } \eta^*(x) = \mathbb{P}(Y = 1 | X = x) \quad g^*(x) = \begin{cases} 1 & \text{if } \eta^*(x) \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{here : If } x, \eta^*(x) = 1/3 < 1/2 \text{ so } g^*(x) = 0.$$

$$\text{Bayes error } \mathbb{P}(Y \neq g^*(X)) = \mathbb{P}(Y = 1) = \frac{1}{3}$$

$$\begin{aligned} \text{in general } \mathbb{P}(Y \neq g^*(X)) &= \mathbb{E}[\mathbb{P}(Y \neq g^*(X) | X)] \\ &= \int P(Y=1) P_X(dx) + \int P(Y=0) P_X(dx) = \\ &\eta^*(x) < 1/2 \quad \eta^*(x) > 1/2 \end{aligned}$$

$$\begin{aligned}
 P(Y \neq g^*(X)) &= \int_{\gamma^* \leq \frac{1}{2}} m^*(x) dP_X(x) + \int_{\gamma^* > \frac{1}{2}} (1 - g^*)(x) dP_X(x) \\
 &= \int \min(\gamma^*, 1 - \gamma^*)(x) dP_X(x) \\
 &= \boxed{\mathbb{E}(\min(\gamma^*, 1 - \gamma^*(x)))}
 \end{aligned}$$



$$k=3: \underline{g_{\text{knn}}(x) = 0}$$

$$k=1 \quad \underline{g_{\text{knn}}(x) = \begin{cases} 1 & \{x \geq 0,55\} \\ 0 & \{x < 0,55\} \end{cases}}$$

errors

$$\begin{aligned}
 R(g_{1,\text{knn}}) &= P(Y \neq \underline{g_{\text{knn}}(x)}) \\
 &= P(Y=1, X < 0,55) \\
 &\quad + P(Y=0, X \geq 0,55) \\
 &= \frac{0,55}{3} + \frac{0,45 \times 2}{3} \\
 &= \frac{1,45}{3} = \boxed{0,48... (> \frac{1}{3})}
 \end{aligned}$$

$$R(g_{3,\text{knn}}) : P(Y=1) = \frac{1}{3}$$

$$5] \quad g(z) = \frac{1}{\pi} \int_{-\pi/2}^{\pi/2} \mathbb{1}\{x \geq \frac{1}{2}\}$$

S-1) $\text{if } b < \frac{1}{2}: P(X \in (a, b), Y=1) = 0$

$\text{if } a > \frac{1}{2} = \frac{b-a}{\pi}$

$\text{if } a < \frac{1}{2} \leq b \quad \dots = P(X \in (\frac{1}{2}, b), Y=1)$

S-2 $g^*(z) = \underline{g(z)}$

$$P(Y = g^*(x)) = \min(\gamma, 1-\gamma) = 0$$

S-3 $R(g_{1,n}) = P(Y=1, g(x)=0) = P(X \in (0, \frac{1}{2}, 0.5))$

$R(g_{3,n}) = P(Y=1) = 0.5$

① $X_n \xrightarrow{P} X \quad \forall \varepsilon > 0 \lim_{n \rightarrow \infty} P(\|X_n - X\| > \varepsilon) \rightarrow 0$

↑
Conv in probability

② $X_n \xrightarrow{a.s} X \quad (\forall \omega \in \Omega)$
almost surely
with respect to measure

$(\exists \Omega \in \mathcal{A}) \quad P(\Omega) = 1$
 $(\forall \omega \in \Omega) \quad \|X_n(\omega) - X(\omega)\| \xrightarrow{n \rightarrow \infty} 0$

$\Leftrightarrow P(\lim_{n \rightarrow \infty} \|X_n - X\| = 0) = 1$

③ $X_n \xrightarrow{D} X \quad (\forall \varepsilon > 0) \quad P(\limsup_{n \rightarrow \infty} \{ \|X_n(\omega) - X(\omega)\| > \varepsilon \}) = 0$

$E(\varphi(X_n)) \xrightarrow{IE(e^{ix_n})} E(\varphi(X)) \quad \forall \varphi \in C_b$

Exercice 2 (Key lemma explaining K-nn consistency):

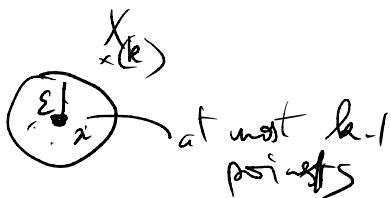
Let P_X be the law of X on \mathbb{R}^d . Let x in the support of P_X , that is for all $\epsilon > 0$, for all ball $B(x, \epsilon)$ of radius ϵ centered at x , $P_X(B(x, \epsilon)) > 0$. Let k_n be a sequence of integers such that $k_n/n \rightarrow 0$. Let $(X_i, i \in \mathbb{N})$ be an iid sequence distributed as X .

- Show that, almost surely, $\|X_{(k_n, n)} - x\| \rightarrow 0$ where $X_{(k_n)}$ is the k_n 'th nearest neighbour of x among (X_1, \dots, X_n) .

Reminder: $x \xrightarrow{\text{a.s.}} X \iff \limsup_{n \rightarrow \infty} P(\limsup_{\omega} \{\omega : |X_n(\omega) - X(\omega)| > \epsilon\}) = 0$

where for a sequence of sets $(A_n)_{n \geq 1}$, $\limsup A_n = \inf_{k \geq 1} \sup_{n \geq k} A_n$

$$\limsup A_n = \inf_{k \geq 1} \sup_{n \geq k} A_n$$



$$\bigcap_{R \geq 1} \bigcup_{n \geq R} A_n$$



here: let $\epsilon > 0$ and $A_\epsilon = \{ \|X_{(k_n, n)} - x\| > \epsilon \}$

$$A_\epsilon \subset \bigcap_{i=1}^n \{X_i \in B(x, \epsilon)\} \leq k_n^{-1} \leq \frac{1}{k_n}$$

$$\subset \left\{ \frac{1}{n} \sum_{i=1}^n \{X_i \in B(x, \epsilon)\} \right\} \leq \frac{k_n}{n}$$

or $\text{L}(\infty) \Rightarrow$ avec proba 1 (over Ω_1 a.s. $\mathbb{P}(\omega \in \Omega_1)$)

$$\frac{1}{n} \sum_{i=1}^n \{X_i \in B(x, \epsilon)\} \rightarrow \underline{P}_\epsilon = \underline{P}(X \in B(x, \epsilon))$$

Now, $\forall n \in \Omega_1$, $\exists n_0 \text{ s.t. } n > n_0 \Rightarrow \underline{P}_{\epsilon/2} > 0$

$$\frac{1}{n} \sum_{i=1}^n \{X_i(\omega) \in B(x, \epsilon)\} \geq \underline{P}_\epsilon > \frac{\epsilon}{2}$$

$\exists m \text{ s.t. } n > m \Rightarrow \frac{k_m}{m} < \frac{\epsilon}{2}$

$$n_2 = \max(n_0, m) \Rightarrow \omega \notin \bigcup_{n \geq n_2} A_n \quad \left(\frac{k_{n_0}}{n_0} < \frac{\epsilon}{2} \right)$$

n chosen such that

$\Rightarrow \omega \notin \limsup A_n$

whence $\limsup A_n \subset \Omega_1^c$

$\underline{P}(\limsup A_n) = 0$ of probability 0.

Exercice 3 (Nadaraya-Watson Regression):

Consider a regression problem for a random pair $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ with target Y and covariate X . We assume that $\mathbb{E}[Y^2] < \infty$. The goal is to approach $m(x) = \mathbb{E}[Y|X = x]$. We assume that the pair (X, Y) has a density $f(x, y)$ with respect to the Lebesgue measure on \mathbb{R}^{d+1} .

1. Write $m(x)$ as an integral involving the f and the marginal density f_X of X .
2. Given kernels K_x, K_y of order 1 for density estimation of X and Y , (thus $\int uK(u)du = 0$ and $\int K(u)du = 1$) define the product kernel density estimate as

$$\hat{f}(x, y) = \frac{1}{nh^2} \sum_{i=1}^n K_x\left(\frac{X_i - x}{h}\right) K_y\left(\frac{Y_i - y}{k}\right).$$

— 1 —

$$\begin{aligned} 1. \quad m(x) &= \mathbb{E}(Y|X=x) = \int_{\mathbb{R}} y f_{Y|X}(y|x) dy \\ &= \int_{\mathbb{R}} y \frac{f(x, y)}{f(x)} dy \end{aligned}$$

- 2.
- (a) Recall the expression for the Kernel density estimate \hat{f}_X of f_X based on K_x and a dataset $X_{1:n}$.
 - (b) Show that the Nadaraya-Watson estimate $\hat{m}(x)$ based on K_x is the plug-in estimate of $\mathbb{E}[Y|X = x]$ based on the expression found in Question 1) up to replacing f_X, f with the kernel density estimates \hat{f}_X and \hat{f} .

$$a) \hat{f}_X(x) = \frac{1}{nh} \sum_{i=1}^n K_x\left(\frac{X_i - x}{h}\right)$$

$$b) \text{recall } \hat{m}(x) = \frac{\sum_{i=1}^n Y_i K_x\left(\frac{X_i - x}{h}\right)}{\sum_{i=1}^n K_x\left(\frac{X_i - x}{h}\right)}$$

with Nadaraya-Watson

$$= \frac{\frac{1}{nh} \sum_{i=1}^n Y_i K_x\left(\frac{X_i - x}{h}\right)}{\frac{1}{nh} \sum_{i=1}^n K_x\left(\frac{X_i - x}{h}\right)}$$

let $\tilde{N}(x) = \frac{1}{nh} \sum_{i=1}^n Y_i K_x\left(\frac{X_i - x}{h}\right)$ be the numerator.

On the other hand the plug-in estimate is:

$$\hat{m}_{\text{PI}}(x) = \int_R y \frac{\hat{f}(x, y)}{\hat{f}(x)} dy$$

$$= \frac{\hat{N}(x)}{\frac{1}{nh} \sum_{i=1}^n K\left(\frac{x_i - x}{h}\right) Y_i}$$

, where :

$$\hat{N}(x) = \int_R y \frac{1}{nh^2} \sum_{i=1}^n K_x\left(\frac{x_i - x}{h}\right) K_y\left(\frac{y - \bar{y}}{h}\right) dy$$

$$= \frac{1}{nh^2} \sum_{i=1}^n K_x\left(\frac{x_i - x}{h}\right) \int_R y K_y\left(\frac{y - \bar{y}}{h}\right) dy$$

change variables: $u = y - \bar{y}$

$$= \frac{1}{nh^2} \sum_{i=1}^n K_x\left(\frac{x_i - x}{h}\right) \int_R (y_i - \bar{y} - hu) K_y(u) du \cdot h \quad [dy = h du]$$

$$= \frac{1}{nh} \sum_{i=1}^n K_x\left(\frac{x_i - x}{h}\right) Y_i$$

$$= \tilde{N}(x)$$

whence $\hat{m}_{\text{PI}}(x) = \hat{m}_{\text{NW}}(x)$

□

$$\left(\frac{\hat{N}(x)}{\hat{f}(x)} \right)$$

$$\left(\frac{\tilde{N}(x)}{\hat{f}(x)} \right)$$

$$\hat{m}(x) = \frac{\sum Y_i K\left(\frac{x_i - x}{h}\right)}{\sum K\left(\frac{x_i - x}{h}\right)}$$

