

Modèles linéaires

* la variable réponse : Y

* les variables explicatives : $X^{(1)}, \dots, X^{(P)}$

question : existe-t'il une fonction f telle que :

$$Y \approx f(X^{(1)}, \dots, X^{(P)})$$

Pour déterminer f , on utilise le risque et un problème de minimisation

On cherche f tq

$$\mathbb{E} \left[\left(Y - f(x^{(1)}, \dots, x^{(P)}) \right)^2 \right] = \min_g \mathbb{E} \left[\left(Y - g(x^{(1)}, \dots, x^{(P)}) \right)^2 \right]$$

Plutôt que de faire la minimisation sur l'ensemble de f et g possibles, on se limite aux fonctions linéaires

$$\mathcal{T}^P = \left\{ g: \mathbb{R}^P \rightarrow \mathbb{R} \text{ tq } g(x_1, \dots, x_P) = \beta_0 + \sum_{k=1}^P \beta_k x_k \right\}$$

pourquoi cette famille.

Car dans le cadre où $Y, X^{(1)}, \dots, X^{(p)}$ sont des variables quantitatives, en réalité, on connaît de façon théorique la meilleure fonction f il s'agit de

$$f(X^{(1)}, \dots, X^{(p)}) = E \left[Y \mid X^{(1)}, \dots, X^{(p)} \right]$$

et si l'on est dans le cadre gaussien, à savoir que le modèle mathématique

est $Y = f(X^{(1)}, \dots, X^{(p)}) + \varepsilon$ avec ε une variable gaussienne, alors l'espérance conditionnelle est une f^{ct} linéaire.

P^b cette espérance conditionnelle est bien souvent incalculable.

On remplace le risque par le risque empirique car à la base, on dispose d'observations.

échantillon d'apprentissage: $\mathcal{L} = \left\{ (x_i^{(1)}, \dots, x_i^{(P)}, y_i) : i \in [1, n] \right\}$

- y_i : une observation de la variable aléatoire Y_i .

- $x_i^{(k)}$ une observation de la variable $X_i^{(k)}$ que nous supposons non aléatoire

le modèle s'écrit: $\forall i \in [1, n], Y_i = \beta_0 + \sum_{k=1}^P \beta_k x_i^{(k)} + \epsilon_i$ la variable de bruit

hypothèses classiques:

$$\forall i \in [1, n], \mathbb{E}[\varepsilon_i] = 0$$

$$\forall i \in [1, n], \text{Var}[\varepsilon_i] = \sigma^2 \quad (\text{homocédasticité})$$

$$\forall i \neq j, \text{Cov}(\varepsilon_i, \varepsilon_j) = 0$$

bien souvent, on ajoute l'hypothèse de normalité du bruit:

$$U = \underbrace{\begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\text{vecteur aléatoire}} \sim \mathcal{P}\left(0_n, \underbrace{\sigma^2 \mathbf{I}_n}_{\text{matrice de variance}}\right) \quad \text{vecteur gaussien}$$

def vecteur gaussien

Soit $T = \begin{pmatrix} T_1 \\ \vdots \\ T_n \end{pmatrix}$ un vecteur aléatoire (T_1, \dots, T_n sont des variables aléatoires)

T est un vecteur gaussien si et seulement si :

$\forall (\lambda_1, \dots, \lambda_n) \in \mathbb{R}^n, \sum_{k=1}^n \lambda_k T_k$ est une variable gaussienne

Rq: Automatiquement $\forall i \in \{1, \dots, n\}, T_i$ est une variable gaussienne

def: matrice de variance d'un vecteur aléatoire T :

$$V[T] = E \left[(T - E[T]) \cdot (T - E[T])' \right] : \text{matrice } n \times n \text{ positive et symétrique}$$

def: espérance d'un vecteur aléatoire T

$$E[T] = \begin{pmatrix} E[T_1] \\ \vdots \\ E[T_n] \end{pmatrix}$$

Rq: $V[T] = \left(\text{cov}(T_i, T_j) \right)_{i,j \in [1, n]}$

la théorie des modèles linéaires est plus vaste

type variables explicatives type de la variable réponse	quantitatives	qualitatives	mélange
quantitative	→ modèle linéaire par moindres carrés (OLS)	Analyse de la variance (ANOVA)	Analyse de la covariance (ANCOVA)
qualitative	Régression Logistique		

$$\triangle! \quad Y_i = \beta_0 + \sum_{k=1}^p \beta_k x_i^{(k)} + \varepsilon_i$$

$$\text{on a } E[Y_i] = \beta_0 + \sum_{k=1}^p \beta_k x_i^{(k)}$$

les Y_i ne sont pas identiquement distribuées car déjà
pas toutes la même espérance!

donner

Variablen

Individuen

y_1

$x_1^{(1)}$

...

$x_1^{(p)}$

y_2

y_n

$x_n^{(1)}$

$x_n^{(p)}$

Rappel:

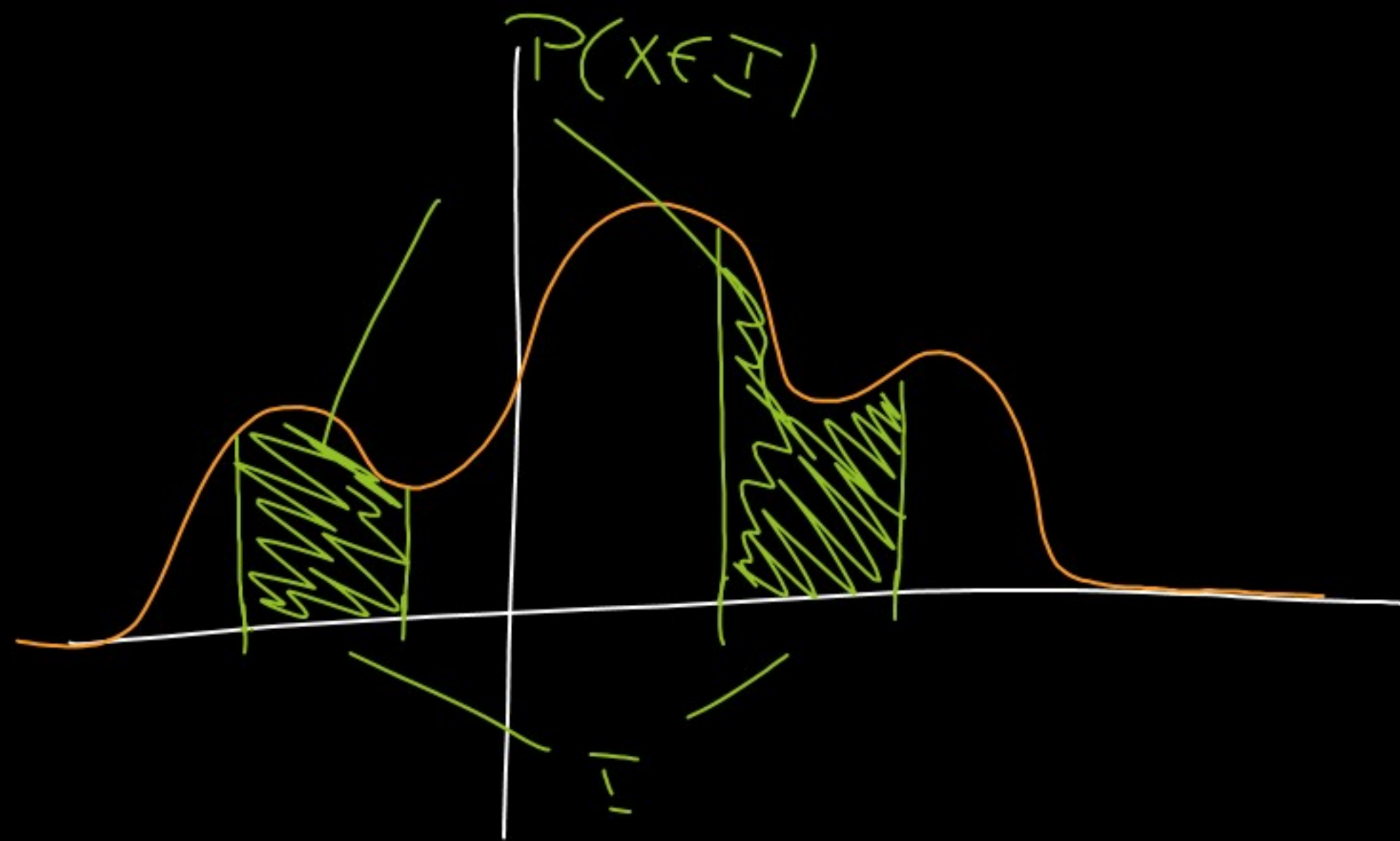
fonction de densité f :

- $f(x) \geq 0 \quad \forall x \in \mathbb{R}$

- $\int_{\mathbb{R}} f(x) dx = 1$

\Rightarrow si X est une variable aléatoire admettant f pour f^d de densité

$$\forall I \subset \mathbb{R}, \quad P(X \in I) = \int_I f(x) dx$$

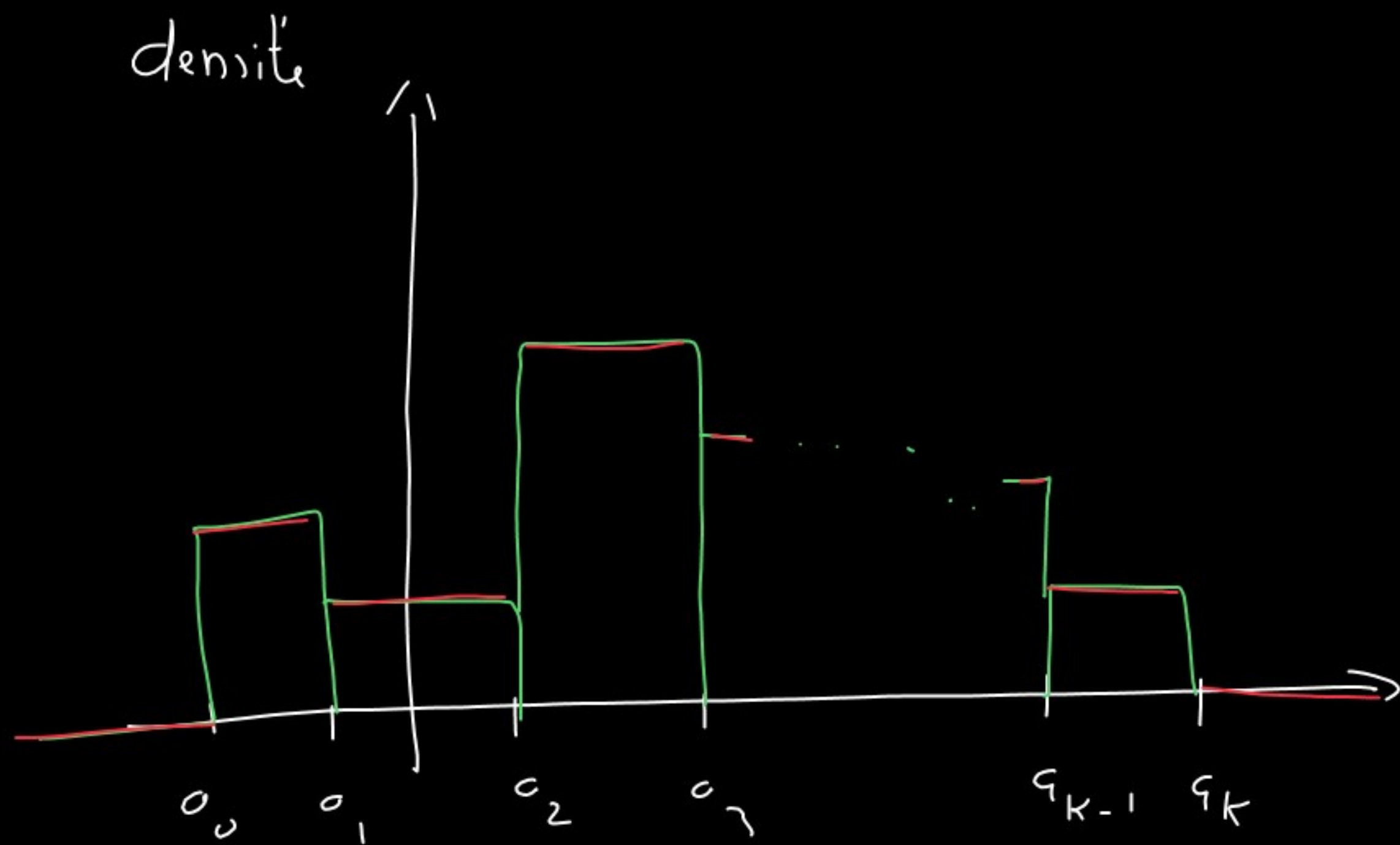


histogramme :

classes	effectif	densité
$[a_0; a_1[$	n_1	d_1
$[a_1; a_2[$	n_2	d_2
$[a_2; a_3[$	n_3	d_3
\vdots	\vdots	\vdots
$[a_{K-1}; a_K[$	n_K	d_K
	<hr/>	
	n	

n_i : n° d'observations qui
sont dans $[a_{i-1}; a_i[$

$$d_i = \frac{n_i}{n(a_i - a_{i-1})}$$



$$\int_{\mathbb{R}} \hat{f}(t) dt = \sum_{k=1}^K d_k (a_k - a_{k-1})$$

$$= \sum_{k=1}^K \frac{n_k}{n} = 1$$

$$\hat{f}(t) = \begin{cases} d_k & \text{pour } t \in [a_{k-1}, a_k] \\ 0 & \text{sinon} \end{cases} \rightarrow f^{\text{d-}} \text{ densité}$$

Estimation de la densité par la méthode des noyaux

def: Noyau (Kernel)

fonction réelle K telle que :

$$\forall x \in \mathbb{R}, K(x) \geq 0$$

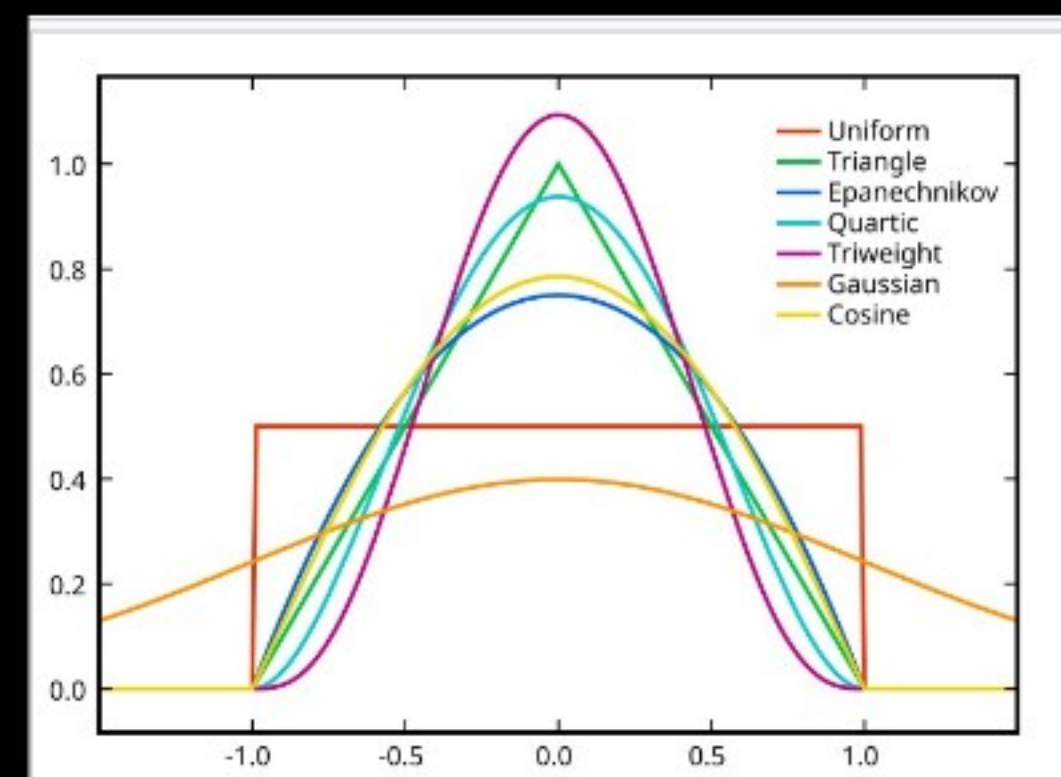
$$\forall x \in \mathbb{R}, K(-x) = K(x)$$

$$\int_{\mathbb{R}} K(x) dx = 1$$

Exemples :

$$\bullet K(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}}$$

$$\bullet K(x) = \frac{1}{2} \mathbb{1}_{[-1;1]}(x)$$



Triangle [modifier | modifier le code]

La forme du noyau est une [fonction triangulaire](#) :

$$K(u) = (1 - |u|) \mathbf{1}_{(|u| \leq 1)}$$

La fonction estimée sera alors linéaire par morceaux.

Epanechnikov [modifier | modifier le code]

On parle aussi de noyau « parabolique ». Il porte le nom de V.A. Epanechnikov, 1969¹ :

$$K(u) = \frac{3}{4}(1 - u^2) \mathbf{1}_{(|u| \leq 1)}$$

Ce noyau permet d'avoir l'estimateur le plus efficace pour la densité.

Quartique [modifier | modifier le code]

$$K(u) = \frac{15}{16}(1 - u^2)^2 \mathbf{1}_{(|u| \leq 1)}$$

Cubique [modifier | modifier le code]

$$K(u) = \frac{35}{32}(1 - u^2)^3 \mathbf{1}_{(|u| \leq 1)}$$

Gaussien [modifier | modifier le code]

$$K(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}u^2}$$

Circulaire [modifier | modifier le code]

autres exemples de
noyaux

Estimateur par noyau

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{k=1}^n K\left(\frac{x - x_i}{h}\right)$$

tuning parameter : bandwidth
(fenêtre)

↳ à calibrer

où les x_i sont les observations de la variable dont on cherche à estimer la densité

Rappel: loi du χ^2 (Chi-Deux)

Soit T_1, \dots, T_d des variables iid $\mathcal{P}(0,1)$

On pose :

$$C = \sum_{k=1}^d T_k^2$$

$$C \sim \chi^2(d)$$

C suit une loi du χ^2 à d degrés de liberté

Test d'adéquation du χ^2 :

→ on crée une distance entre les effectifs observés

et les effectifs théoriques, les effectifs théoriques étant

calculés sous l'hypothèse H_0 qui précise totalement la loi inférée pour la variable.

On peut alors définir la statistique du χ^2 :

$$T = \sum_{j=1}^J \frac{(N\hat{p}_j - Np_j)^2}{Np_j} = \sum_{j=1}^J \frac{(n_j - Np_j)^2}{Np_j} \text{ où } n_j = N\hat{p}_j = \sum_{i=1}^N 1_{y_i=v_j}$$

Résultat:

Sous H_0 , $T \underset{\mathcal{L}}{\sim} \chi^2(J-1)$

La règle de décision est :

si $T > c \Rightarrow$ on rejette H_0

si $T < c \Rightarrow$ on ne rejette pas H_0

avec $c = c_{1-\alpha; J-1}$: le quantile d'ordre $1-\alpha$ pour $\chi^2(J-1)$

$$(P(\chi^2(J-1) > c_{1-\alpha; J-1})) \leq \alpha$$

α est ce que l'on appelle le niveau de test.

niveau d'un test

réalité décision \	H_0	H_1
H_0	OK	erreur
H_1	erreur	OK

P_{H_1} (ne pas rejeter H_0)

" 1 - puissance

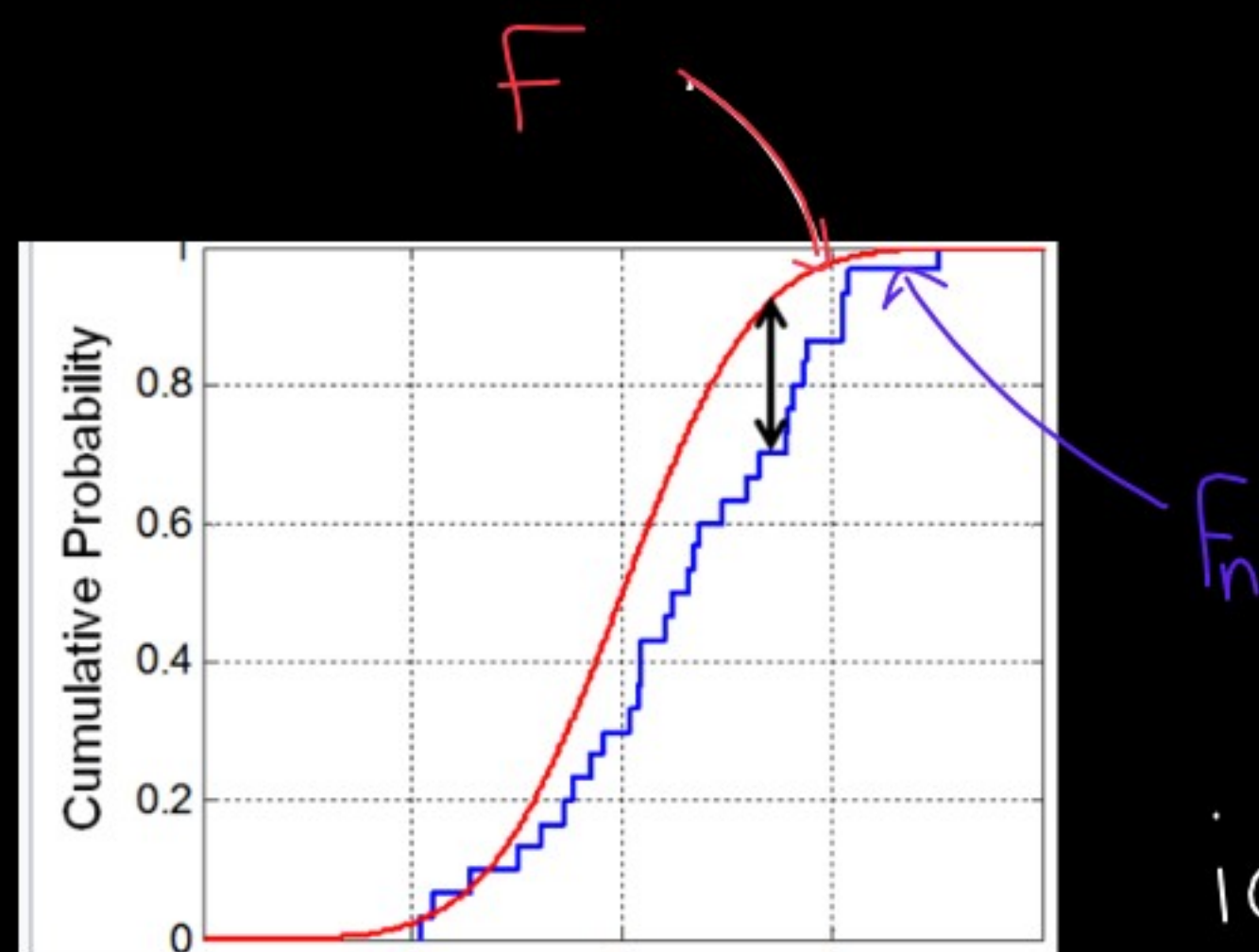
↑ on minimise cette erreur : $P_{H_0}(\text{décider } H_1) \leq \alpha$

test de Kolmogorov

On dispose d'observations x_1, \dots, x_n

H_0 : les observations sont associées à une variable aléatoire de f^d de répartition F

H_1 : ce n'est pas le cas



On introduit la f^d de répartition empirique

$$F_n(t) = \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{X_k \leq t}$$

idée: créer une distance entre F et F_n

la distance est:

$$K_n = \sup_{t \in \mathbb{R}} |F(t) - F_n(t)|$$

Sous H_0 ,

$$\sqrt{n} K_n \rightsquigarrow K \quad \nwarrow \text{loi de Kolmogorov}$$

Région de rejet:

$$\left\{ \sqrt{n} K_n > \alpha \right\} \text{ avec } n \text{ tq } P(K > \alpha) < \alpha$$

En pratique, on ne calcule pas τ mais la p-value.

Pour le test de Kolmogorov, la p-value est donnée par:

$$p\text{-value} = P\left(K > \sqrt{n} \underbrace{K_{n,obs}}\right)$$

↓
la valeur calculée de K_n
Sur le jeu de données

Règle de décision:

- si $p\text{-value} < \alpha \Rightarrow$ on décide H_1
- si $p\text{-value} > \alpha \Rightarrow$ on ne rejette pas H_0