

(S5) Estimation "plus lisse" d'une densité de probabilité  
(plus lisse que l'histogramme)

Cadre une v.a.r  $Y$  et ses réalisations  $y \in \mathbb{R}$

indices: 1 2 ...  $l_2$  ...  $N$

observations:  $y_1 y_2 \dots y_{l_2} \dots y_N$

ex  
 $N=16$

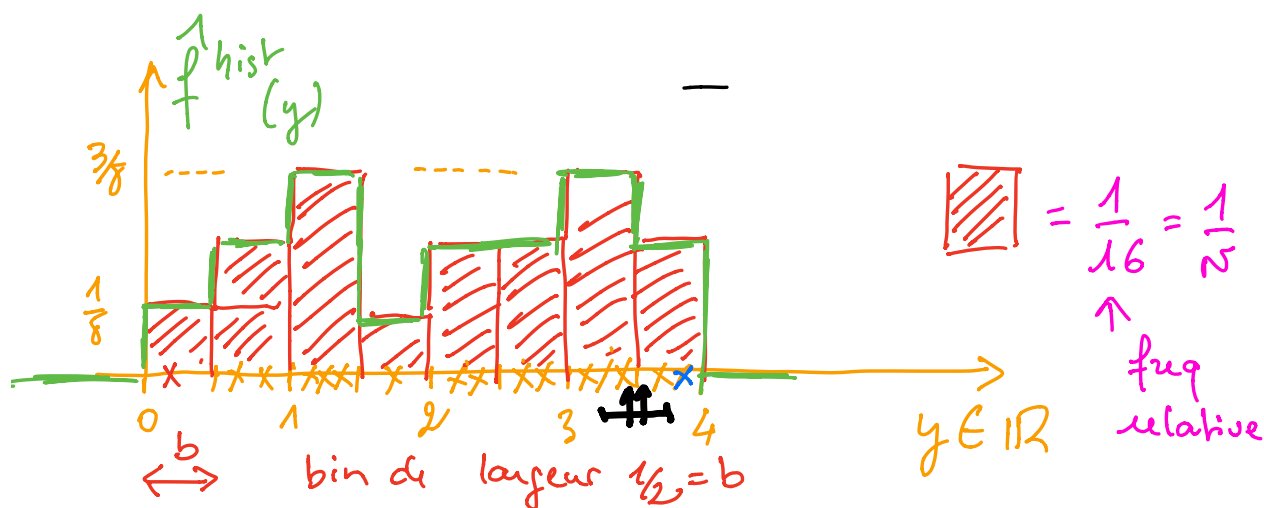
1 Réalisation  
d'un 16  
échantillon

$x$   $x_{\min}$  ...  $x_{\max}$  ...  $x$

ex  $\in [0, 4[$

choix des 8 intervalles / "bins" de largeur  $b = 1/2$

Classe	$[0; 0,5[$	$[0,5; 1[$	$[1; 1,5[$	$[1,5; 2[$	...	$[3,5; 4[$
Effectif	1	2	3	1	...	2
Freq Relative	$1/16$	$2/16$	$3/16$	$1/16$		$2/16$



Une formulation plus "mathématique"  
du calcul de l'histogramme

$$\hat{f}^{\text{hist}}(y) = \begin{matrix} \text{densité en } y \\ \text{vraisemblance de } y \end{matrix}$$

$$\int_{\mathbb{R}} \hat{f}^{\text{hist}}(y) dy = \underbrace{b \times \frac{1}{8}}_{\frac{1}{16}} + \underbrace{b \times \frac{2}{8}}_{\frac{2}{16}} + b \times \frac{3}{8} + \dots$$

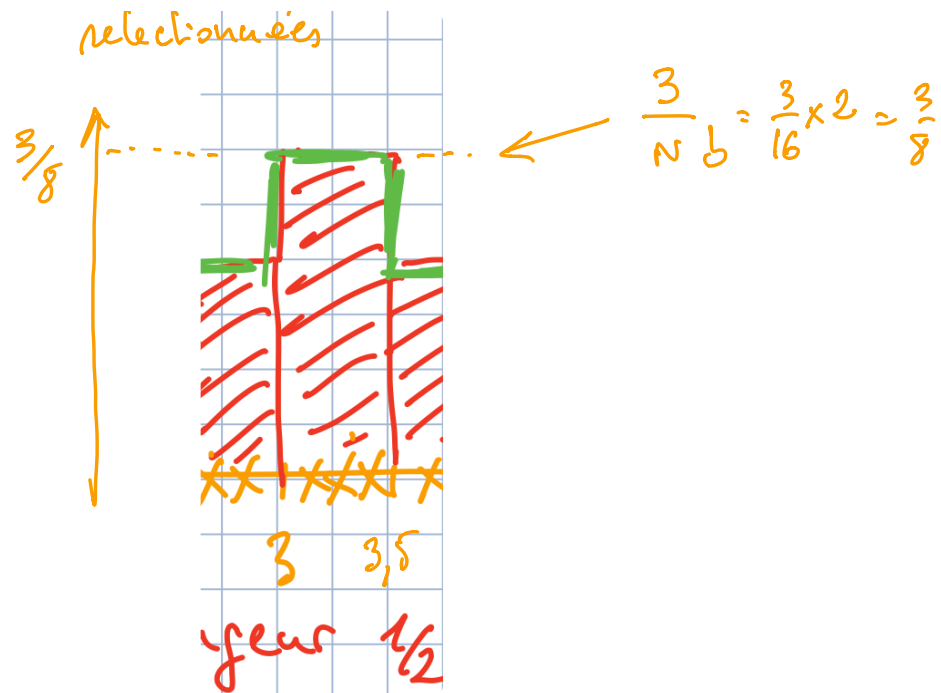
$$= 1$$

Une formule qui traduit le tracé  
de l'histogramme

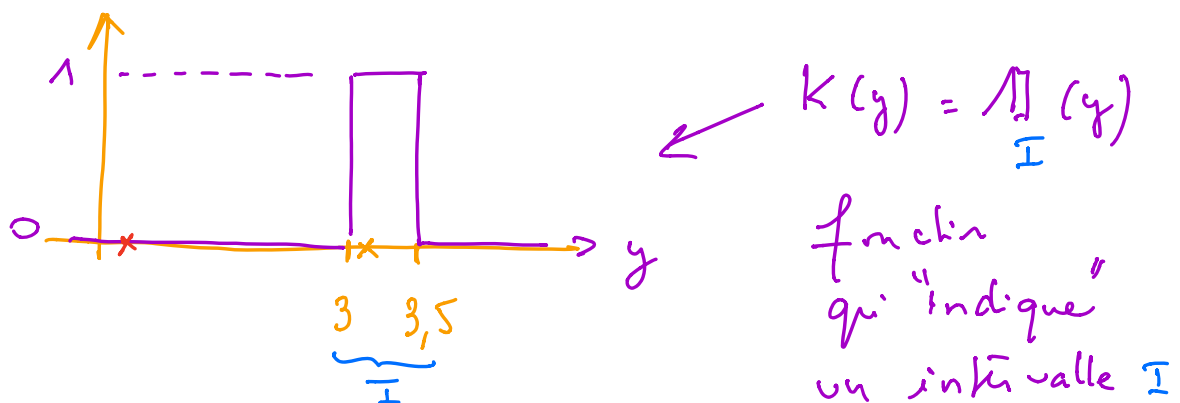
"selectionner les réalisations appartenant  
au bin, les compter, diviser par  
la taille de l'échantillon ( $N$ ),  
diviser par la largeur  $b$  du bin"

exemple bin  $\overset{\leftarrow b \rightarrow}{[3; 3,5[}$

3 réalisations  $\quad \times \times \times$



Fonction indicatrice d'un intervalle  $I$   
d'un bin  $I$



$$1_I(y) = \begin{cases} 1 & \text{si } y \in I \\ 0 & \text{sinon} \end{cases}$$

Compter les

$\sum_{i=1}^N 1_I(y_i) = 3 !$

16 dans notre exemple

réalisations

dans  $I$

Diviser par  $N$

(Pour trouver la freq relative)

Diviser par  $b$

(Pour respecter les unités de surface  
et assurer ainsi  $\int f = 1$ )

$$\frac{1}{Nb} \sum_{i=1}^N \mathbb{1}_I(y_i) = \frac{3}{16 \frac{1}{2}} = \frac{3 \times 2}{16} = \frac{3}{8}$$

$$\hat{f}^{high}(y) = \frac{1}{N} \frac{\# \text{ de } y_i \text{ dans le bin que } y}{\underbrace{\text{largeur du bin}}_b}$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{1}{b} \mathbb{1}_{\text{Bin}(y)}(y_i)$$

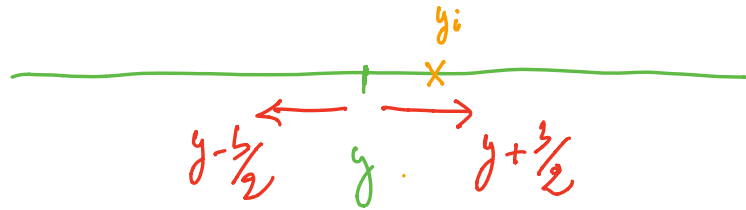
Mais alors on peut aussi faire glisser la  
fenêtre de comptage (au lieu de fixer les bins)

$$\hat{f}^{\text{fenêtre glissante}}(y) = \frac{1}{N} \sum_{i=1}^N \frac{1}{b} \mathbb{1}_{[y-\frac{b}{2}; y+\frac{b}{2}]}(y_i)$$

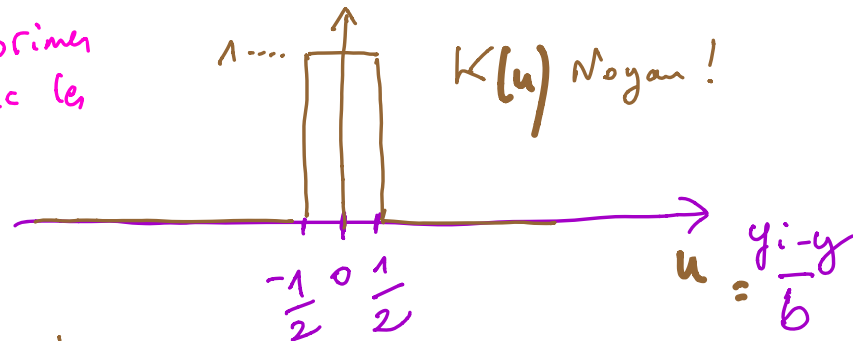
inconvenient

↗ c'est l'intervalle qui

depend de  $y$



Idee : réexprimer la fenêtre avec les écarts



$$K\left(\frac{y_i - y}{b}\right) = 1 \quad (\Leftrightarrow) \quad -\frac{1}{2} \leq \frac{y_i - y}{b} \leq \frac{1}{2}$$

$$= \mathbb{1}_{\left[\frac{y - b/2}{b}, \frac{y + b/2}{b}\right]} \quad -\frac{b}{2} \leq y_i - y \leq \frac{b}{2}$$

$$\left[-\frac{1}{2}, \frac{1}{2}\right]$$

$$y - \frac{b}{2} \leq y_i \leq y + \frac{b}{2}$$

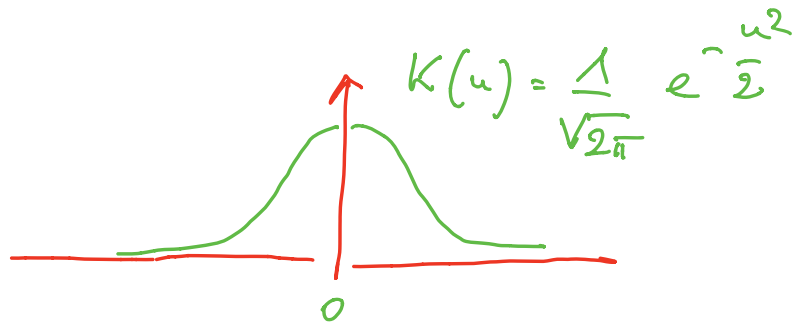
$$\hat{f}_{\text{fenêtre glissante}}(y) = \frac{1}{N} \sum_{i=1}^N \frac{1}{b} \mathbb{1}_{[y - b/2; y + b/2]}(y_i)$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{1}{b} K\left(\frac{y_i - y}{b}\right)$$

Symétriques centrés en les  $y_i$   $\sum$  de  $N$  "noyaux"

Mais alors d'autres noyaux  $K$  (Kernel) sont possibles

$$\hat{f}^{\text{noyau}}(y) = \frac{1}{N} \sum_{i=1}^N \frac{1}{b} K\left(\frac{y_i - y}{b}\right)$$



Noyau gaussien !

Effet de  $b$  !!

Regardez une vidéo sur "Kernel Density Estimation" bien choisir

## Choice of kernel

The *triangular* kernel (or *linear* kernel) is given by

$$f(x) \propto \max(1 - |x|, 0).$$

