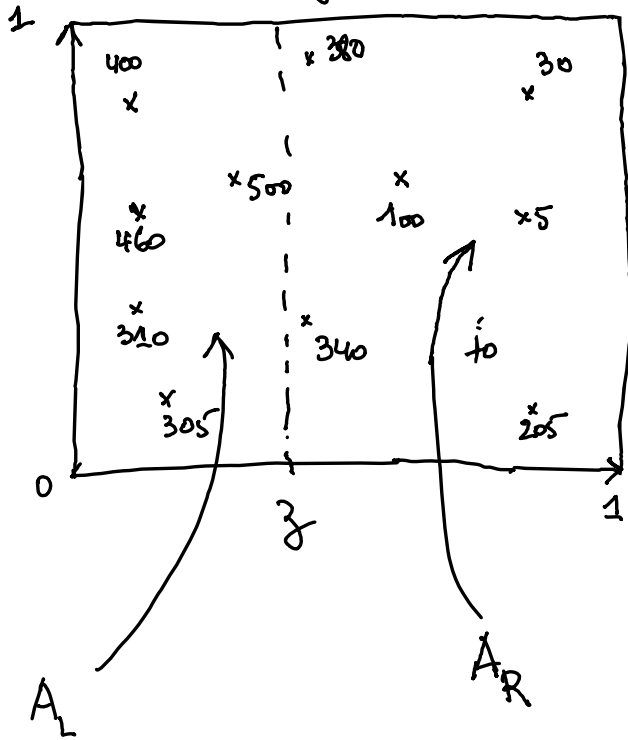


# Arbre de régression : (Construction)

Critère :



$$L_n(i, z) = \frac{1}{N_n(A)} \sum_{i=1}^n (y_i - \bar{y}_{A_L} \mathbb{1}_{\{X_i^{(j)} < z\}} - \bar{y}_{A_R} \mathbb{1}_{\{X_i^{(j)} \geq z\}})^2$$

- $A_L = \{x \in A : x^{(j)} < z\}$
- $A_R = \{x \in A : x^{(j)} \geq z\}$
- $\bar{y}_A$  est la moyenne de  $y_i$  qui se trouvent dans  $A$ .
- $N_n(A)$  le nombre d'observations dans  $A$ .

Le jeu de données

$$D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

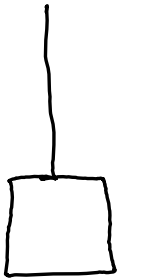
$$x \in \mathbb{R}^d \quad \text{et} \quad y \in \mathbb{R} \quad (\text{régression}).$$

Split d'un poids  $A$ : pour toute variable  $j$

Calculer  $L_A(\bar{j}, z_j)$  comme suit

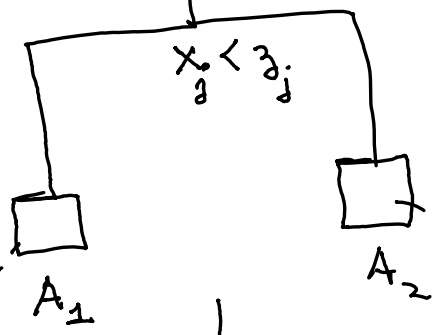
$j$ :	soit $z_j$	$L_A(j, z_j)$
1	$z_1$	$L_A(1, z_1)$
2	$z_2$	$L_A(1, z_2)$
$\vdots$	$\vdots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$
$d$	$z_d$	$L_A(d, z_d)$

on retient le split qui correspond à  
arg min  $L_A(\bar{j}, z_j)$   
 $\bar{j} \in \{1, \dots, d\}$



$A = D_n$

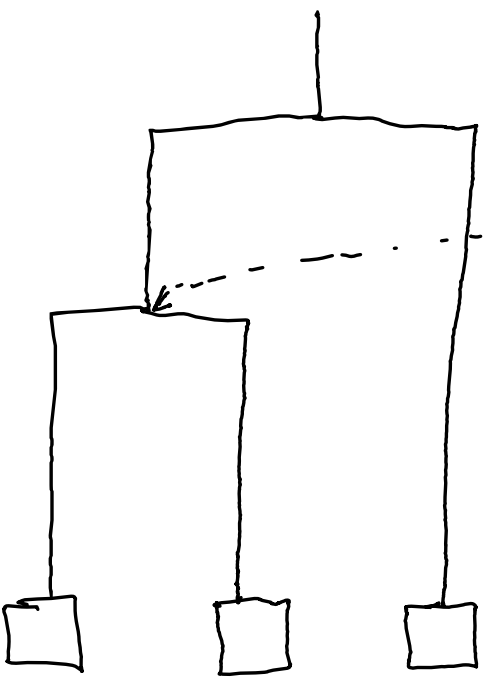
split  
Iteratum ①



$A_1$

$A_2$

évolution de la piste



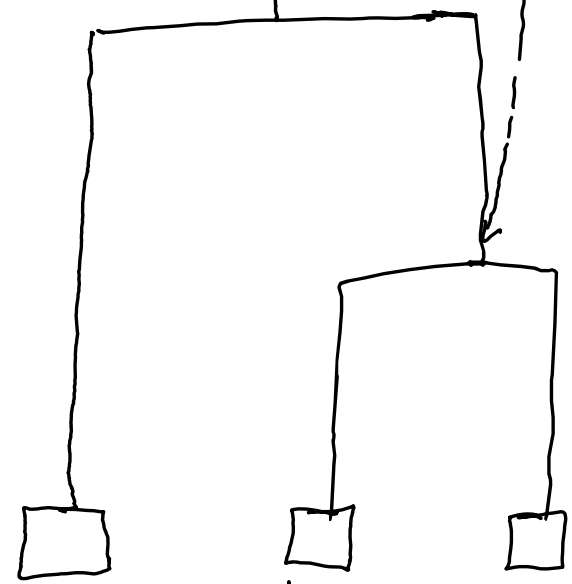
$A_{1,1}$

$A_{1,2}$

$A_2$

Iteratum ②

on garde  
la meilleure évolution  
de la piste !!

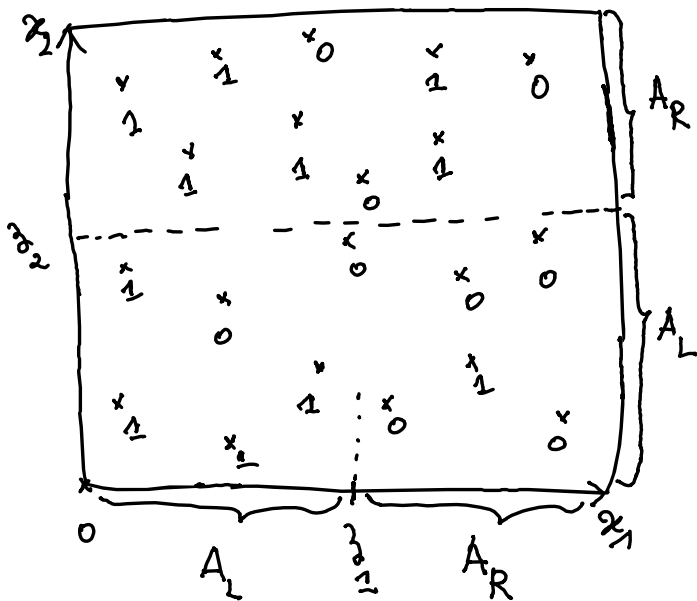


$A_1$

$A_{2,1}$

$A_{2,2}$

En classification : le jeu de données  $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$   
 $x \in \mathbb{R}^d$  et  $y \in \{0, 1\}$  (classification).



Indice ou mesure d'impureté de Gini:

$$\hat{\phi}_A(j, k) = \frac{\sum_{i=1}^n \mathbb{1}\{y_i = k\} \times \mathbb{1}\{x_i \in A\}}{\sum_{i=1}^n \mathbb{1}\{x_i \in A\}}$$

Condition sur  $x_i^{(j)}$

$$L_A(j, \mathcal{Z}) = \underbrace{\left(1 - \sum_{k=0}^2 \left[ \hat{\phi}_{A_L}(j, k) \right]^2 \right)}_{\gamma_{A_L}} + \underbrace{\left(1 - \sum_{k=0}^2 \left[ \hat{\phi}_{A_R}(j, k) \right]^2 \right)}_{\gamma_{A_R}}$$

$A = A_L \cup A_R$

→ construction de l'arbre identique à la régression en adaptant l'erreur. Impureté au lieu de la perte quadratique moyenne.

Élagage d'un arbre profond:

Proposition: Le paramètre de complexité  $\alpha$  (cp dans  $\mathbb{R}$ )

$$\alpha = \frac{L(t) - L(T_t)}{|\tilde{T}| - 1}$$

$T$ : un arbre  
 $\tilde{T}$ : les feuilles de l'arbre  $T$ .

Rappel:

$$C_\alpha(T) = \underbrace{L(T)}_{\substack{\text{perte} \\ \text{ou erreur}}} + \alpha \underbrace{|\tilde{T}|}_{\text{pénalité}}$$

$T_t$  : est la branche de l'arbre qui part du nœud "t".

pour un nœud  $t \in T$ , nous avons:

$$C_\alpha(t) = L(t) + \alpha \quad \text{car un seul nœud terminal "t".}$$

De manière similaire, pour toute branche  $T_t \subset T$ , nous avons

$$C_\alpha(T_t) = L(T_t) + \alpha |\tilde{T}_t|.$$

$$\text{Lorsque } \alpha = 0, \quad C_0(t) = L(t) > L(T_t) = C_0(T_t).$$

Cette inégalité est vraie parce que la première étape d'élagage - consiste à retirer tous les nœuds terminaux tels que  $L(t) = L(t_L) + L(t_R)$ .

Pour les nœuds restants, nous avons:  $L(t) > L(t_L) + L(t_R)$ .

Si on augmente  $\alpha$ ,  $C_\alpha(T_t)$  va augmenter plus vite que  $C_\alpha(t)$  parce que  $|\tilde{T}_t| > 1$ .

Autrement dit, pour un certain  $\alpha$ , nous allons avoir  $C_\alpha(\tilde{T}_t) = C_\alpha(t)$  - (1)

Donc, on obtient  $L(t) + \alpha = L(\tilde{T}_t) + \alpha |\tilde{T}_t|$

$$\Rightarrow (|\tilde{T}_t| - 1) \alpha = L(t) - L(\tilde{T}_t)$$

$$\Rightarrow \alpha = \frac{L(t) - L(\tilde{T}_t)}{|\tilde{T}_t| - 1}$$

Ensuite, le procédé devient itératif pour obtenir une

suite  $\alpha_1 < \alpha_2 < \alpha_3 < \dots < \dots$