

## I. Introduction

- $\ell: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  local loss function

Gm veut résoudre :  $\arg \min_f \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(y, f(x))]$

- Classification:

(0/1) prediction loss  $\rightarrow$  best classifier is Bayes

$$f_{\text{Bayes}}(x) = \arg \max_c \mathbb{P}(Y = c | x)$$

- Regression:

$\ell_2$  loss:  $(y - f(x))^2 \rightarrow$  best solution in regression is

$$f_{\text{reg}}(x) = \mathbb{E}(Y | x)$$

- Min of empirical risk:  $\arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) + \lambda \Omega(f)$

- Statistical learning: entrée:  $\mathcal{H}, \ell, \Omega, S_n$

sortie:  $f_n = \mathcal{A}(S_n, \mathcal{H}, \ell, \Omega)$

learning algo  $\uparrow$  train data  $\uparrow$  class of  $f$   $\uparrow$  loss  $\uparrow$  complexity measure

Parametric VS nonParametric

$\hookrightarrow$  nb of param  $\uparrow$  with size of training

Tree-based:  
 $\checkmark$  transparency, pre-training, interprétabilité  
 $\times$  dépend bcp du nb de training data

## II. Decision and regression tree

### A. Predictive models

## Decision trees:

- model of the form:  $f(x) = \sum_{\ell=1}^m \mathbb{1}(x \in R_\ell) f_\ell$

→ fonctions en escalier

$$f_{\ell} = \arg \max_c \hat{p}_{\ell c} \quad \text{avec} \quad \hat{p}_{\ell c} = \frac{1}{N_{\ell}} \sum_i \mathbb{1}(y_i = c)$$

- make partition based on local criteria
- orthogonal separators
- leaf = majority vote

## Regression trees:

model of the form :  $g(x) = \sum_{\ell=1}^m 1(x \in \mathcal{R}_\ell) g_\ell$

avec  $y_e = \frac{1}{N_e} \sum y_i$

- leaf = local average

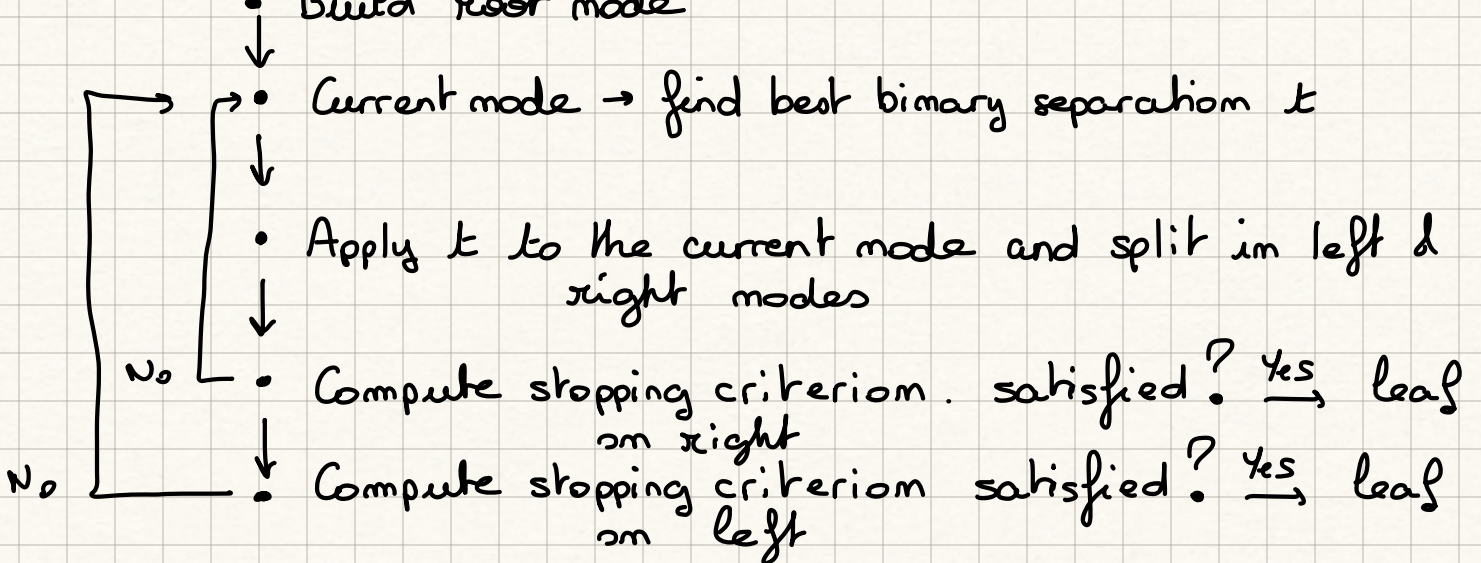
## B. Learning trees

continuous:  $t_{j,s}(x) = \text{sign}(\underbrace{x^j}_{\text{feature}} - \underbrace{s}_{\tau_{\text{split}}})$

### Recursive building algo:

categorical:  $t_{j,v,k}(x) = 1(x^j = v_k^j)$

- Build root mode





## Best split for regression?

2 plans :  $\mathcal{R}_\ell(j, s) = \{x, x^j \leq s\}$   
 $\mathcal{R}_r(j, s) = \{x, x^j > s\}$

→ on veut minimiser  $L((j, s), S)$  cad :

$$\min_{j, s} \left[ \sum_{x_i \in S \cap \mathcal{R}_r} (y_i - \hat{g}_r)^2 + \sum_{x_i \in S \cap \mathcal{R}_\ell} (y_i - \hat{g}_\ell)^2 \right]$$

$\uparrow$  moy des  $y_i$  à d       $\uparrow$  moy des  $y_i$  à g

→ en pratique il suffit de minimiser la variance empirique

$$\min_{j, s} N_r \text{Var}(Y | S \cap \mathcal{R}_r(j, s)) + N_\ell \text{Var}(Y | \mathcal{R}_\ell(j, s))$$

## Best split for classification?

goal = minimize an impurity criterion

$$\min \left[ \frac{N_r}{N} H(S \cap \mathcal{R}_r(j, s)) + \frac{N_\ell}{N} H(S \cap \mathcal{R}_\ell(j, s)) \right]$$

Impurity criteria :  $p_k(S) = \frac{1}{n} \sum_{i=1}^n 1(y_i = k) \rightarrow$  proba de chaque classe

• Cross-entropy :  $H(S) = - \sum_{k=1}^c p_k(S) \log p_k(S)$

• Gini :  $H(S) = \sum_{k=1}^c p_k(S) (1 - p_k(S))$

• Missclassification :  $H(S) = 1 - p_{k(S)}$

Stopping criterion : • maximal depth

• un nœud a 1 mb min de data

↓  
améliorable par cross-validation

Désavantage : très variable en fonction des données  
→ 1 petite variation entraîne un arbre complètement différent

### III. Introduction to ensemble methods

#### How to combine models?



##### Ensemble Method

Same Base Model  
Combinaison linéaire

eg: Bagging, Rd Forest,  
Boosting



##### Committee Machines (Wisdom of the crowd)

Diverse base model  
Combinaison can be  
non-linéaire

↳ learned  
eg: Blending, Stacking



##### Mixture of Experts

Models are experts  
for a certain data  
type

On utilise le bon  
modèle en fonction  
de la data

#### Focus on Ensemble methods

- bootstrap samples (Bagging & Random Forest)
- randomized pred (Rd. F)
- weighted version of sample (Boosting)

### IV. Bagging

Reminder: noise = error by the Bayes model

biais = diff btw minimal error (Bayes error)  
and the avg model

variance = how much  $h_S$  varies from 1 training  
set to another

Algo: ① draw  $T$  bootstrap samples  $\{B_1, \dots, B_T\}$  from  $S$   
uniform sampling with replacement

②  $\forall B_T$ , learn a model  $f_t$

③ Build avg model:  $f_{\text{avg}}(x) = \frac{1}{T} \sum_{t=1}^T f_t(x)$

Useful when learning algo unstable with a high variance  
→ genre les arbres à tout hasard

↳ rarement utilisé seul

↳ pour classifier, un mauvais modèle empire la situation



## V - Random Forests

RF algo : • Boucle de 1 à T :

- ① Bootstrap sample from S
- ② Build randomized DT  $h_{tree}^{(t)}$ 
  - at each node :
    - ① select k features sans remise
    - ② choose the cut

!! do  
not prune

$$\text{out: } H^T(x) = \frac{1}{T} \sum h_{tree}^{(t)}(x)$$

Extra trees algo : • Boucle de 1 à T

- ① Prend Strain en entier
- ② Build randomized DT  $h_{tree}^{(t)}$ 
  - at each node :
    - ① select k features sans remise
    - ② Draw k splits using Pick-a-random-split
      - $a_{max}, a_{min}$  extrema of  $x_i$
      - draw cut point in  $[a_{max}, a_{min}]$
    - ③ Choose the best split des k

$$\text{out: } H^T(x) = \frac{1}{T} \sum h_{tree}^{(t)}(x)$$

RF pros & cons

X - if tree too large → overfitting  
- perte d'interprétabilité

✓ - rapide et marche bien quand ya bcp de features  
- facile à tune

Variable importance : calculé en prenant les variables non tirées.  
Pour tester la feature j il permute dans les data les features j et fait 2 risques, si arbre permuté - bon alors feature importante

$$\hat{f}(x^j) = \frac{1}{n_{tree}} \sum_t R_n(h_t, \bar{S}_n^{t,j}) - R_n(h_t, \bar{S}_m^t)$$

## VI. Boosting

### A. AdaBoost as a Greedy Scheme

Idea: on "boost" en ajoutant un poids aux données en fonction de la loss obtenue jusqu'à now.

Weak classifier: classifier with avg training error  $\leq 0,5$

Algo: Init:  $\omega_1(i) = \frac{1}{n}$  et  $H_0 = 0$

Boucle: ①  $h_t = \operatorname{argmin}_h \sum_1^m \varepsilon_t(h)$  avec  $\varepsilon_t(h) = \mathbb{P}_{i \sim \omega_t}(h(x_i) \neq y_i)$

②  $\alpha_t = \frac{1}{2} \log \frac{1 - \varepsilon_t}{\varepsilon_t}$  avec  $\varepsilon_t = \varepsilon_t(h_t)$

encourage à corriger les trucs mal classifiés  $\leftarrow$  ③  $\omega_{t+1}(i) = \omega_t(i) \frac{e^{-\alpha_t y_i h_t(x_i)}}{Z_{t+1}} = \frac{e^{-y_i H_t(x_i)}}{m \prod_s Z_s}$

④  $H_t = H_{t-1} + \alpha_t h_t$

Out:  $F_T = \operatorname{sigm}(H_T)$

↑  
normalisation  
tg  $\sum \omega_{t+1} = 1$   
 $Z_t = 2 \sqrt{\varepsilon_t (1 - \varepsilon_t)}$

Bounding the training error:  $R_n(F_T) \leq e^{-2 \sum (\frac{1}{2} - \varepsilon_t)^2}$

$R_n(F_T) \leq e^{-2 \gamma^2 T}$  où  $\gamma \leq (\frac{1}{2} - \varepsilon_t)$

### B. Gradient Boosting

Algo: minimizer step deviant: ①  $h_t$  = steepest descent direction in  $\mathcal{H}$

easy if finding the best direction is easy

② choose  $\alpha_t$  minimize  $L(y, H + \alpha h_t)$



Link . AdaBoost :  $\ell(y, h) = e^{-yh}$

• LogitBoost :  $\ell(y, h) = \log(1 + e^{-yh})$

•  $L_2$  Boost :  $\ell(y, h) = (y - h)^2$

•  $L_1$  Boost :  $\ell(y, h) = |y - h|$

• HuberBoost :  $\ell(y, h) = |y - h|^2 \mathbb{1}_{|y - h| < \epsilon} + (2\epsilon |y - h| - \epsilon^2) \mathbb{1}_{|y - h| \geq \epsilon}$