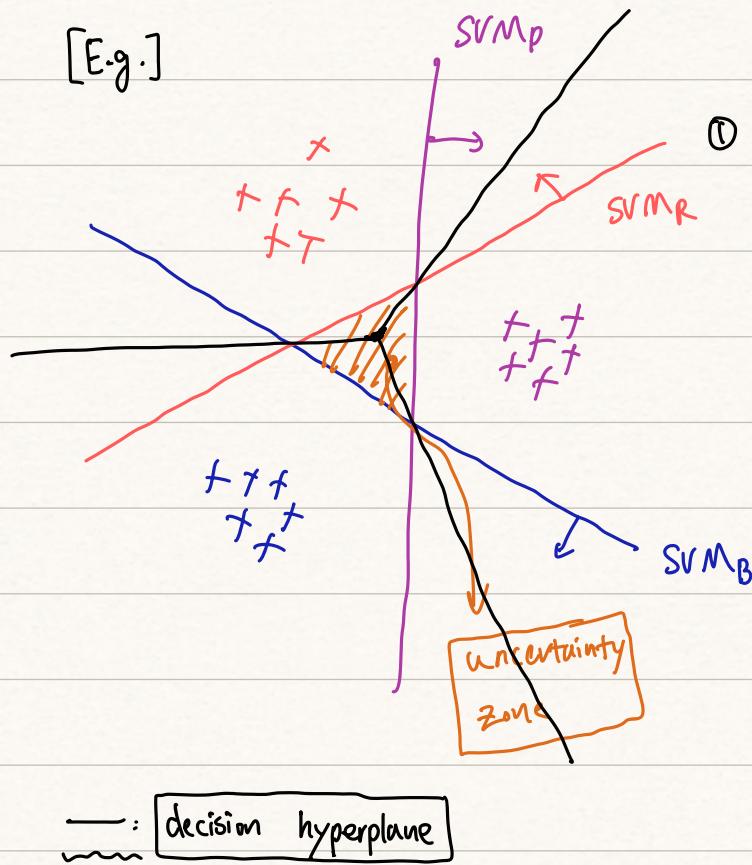


DSA5105 lec4

Recap:

SVM → Multi-class classification

[E.g.]



① 1 vs Rest

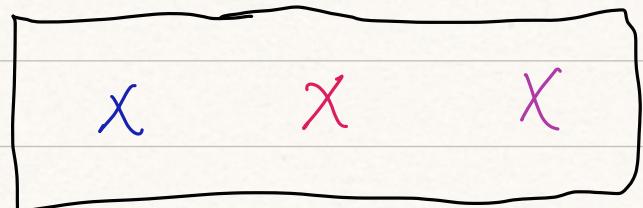
- { SVM_R : Red vs Not Red
- SVM_p : Pur vs Not Pur
- SVM_B : Blue vs Not Blue

IDEA

$$\left\{ \begin{array}{l} w_0^T x + b_0 \\ w_p^T x + b_p \\ w_B^T x + b_B \end{array} \right.$$

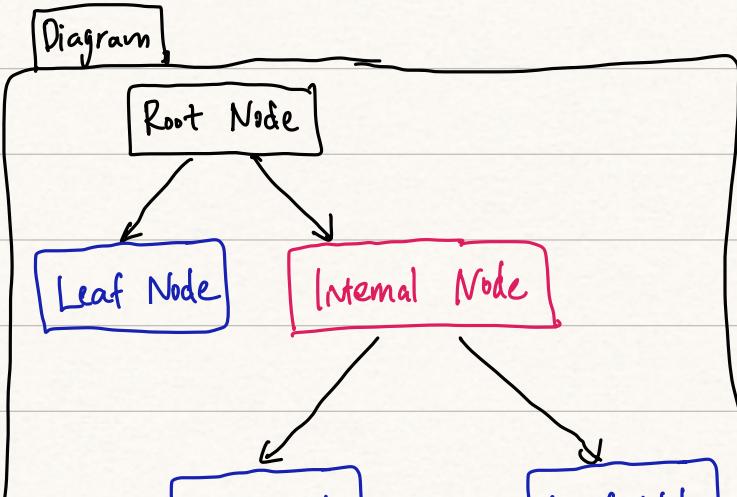
compare these 3 values

↓
Cannot work in such kind of dataset



Hypothesis space $H = \{ f : \text{some kind of function } f \}$

DTs → subsequent questions



Leaf Node

Leaf Node

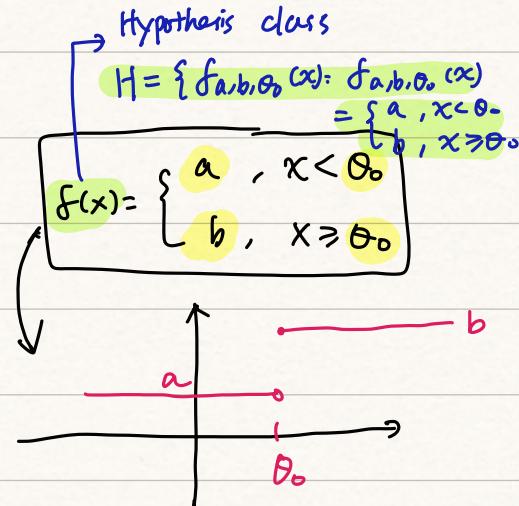
Key Idea:

→ Stratify the input space into non-overlap regions!

[Eg.] → 1-dim example

$$f^*: [0,1] \rightarrow \mathbb{R} \rightsquigarrow \text{Oracle}$$

A depth-1 DT is a piece-wise function:



trainable parameter (a, b, θ_0)

suppose Oracle $f^*(x) = x$ ↼ identical mapping

suppose $x \sim U[0,1]$

$$\text{then } R_{\text{pop}}[f_{a,b,\theta_0}(x)] = \mathbb{E}_{(x,y)} [\text{Loss}(f_{a,b,\theta_0}(x), y)]$$

$$= \frac{1}{2} \int_0^{\theta_0} (a-x)^2 dx + \frac{1}{2} \int_{\theta_0}^1 (x-b)^2 dx$$

consider minimize R_{pop} to achieve $(\hat{a}, \hat{b}, \hat{\theta}_0) = (\frac{1}{2}, \frac{1}{4}, \frac{3}{4})$

$$\textcircled{1} \rightarrow \frac{\partial R_{\text{pop}}}{\partial a} = \frac{(a-\theta_0)^2}{2} - \frac{a^2}{2}$$

$$= \frac{\theta_0(\theta_0-2a)}{2}$$

$$\textcircled{2} \rightarrow \frac{\partial R_{\text{pop}}}{\partial b} = \frac{(1-b)^2}{2} - \frac{(\theta_0-b)^2}{2}$$

$$= \frac{(1-\theta_0)(\theta_0+1-2b)}{2}$$

$$③ \rightarrow \text{first-order optimality condition} \Rightarrow \begin{cases} \frac{\partial R_{\text{pop}}}{\partial a} = 0 \\ \frac{\partial R_{\text{pop}}}{\partial b} = 0 \end{cases} \Rightarrow \begin{cases} \hat{a} = \frac{1}{2}\hat{\theta}_0 \\ \hat{b} = \frac{1}{2}(1 + \hat{\theta}_0) \end{cases}$$

$$④ \rightarrow \text{find } \hat{\theta}_0, \quad \frac{\partial R_{\text{pop}}}{\partial \theta_0} = \frac{1}{2}(a - \theta_0)^2 - \frac{1}{2}(\theta_0 - b)^2$$

$$\text{first-order optimality condition} \Rightarrow \frac{\partial R_{\text{pop}}}{\partial \theta_0} = 0$$

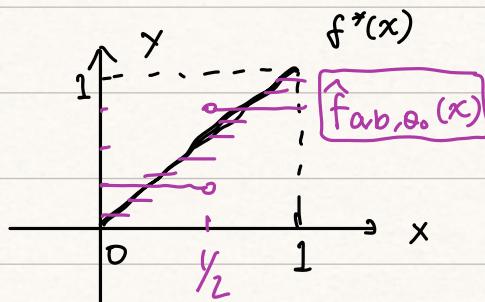
$$\Leftrightarrow (\hat{a} - \hat{\theta}_0)^2 = (\hat{b} - \hat{\theta}_0)^2$$

$$\Leftrightarrow \frac{1}{4}\hat{\theta}_0^2 = \left(\frac{1}{2} - \frac{1}{2}\hat{\theta}_0\right)^2$$

$$\Leftrightarrow \boxed{\hat{\theta}_0 = \frac{1}{2}}$$

$$\Rightarrow \begin{cases} \hat{a} = \frac{1}{4} \\ \hat{b} = \frac{3}{4} \end{cases}$$

That is



use more trees to make the approximation more accurate

Another example:

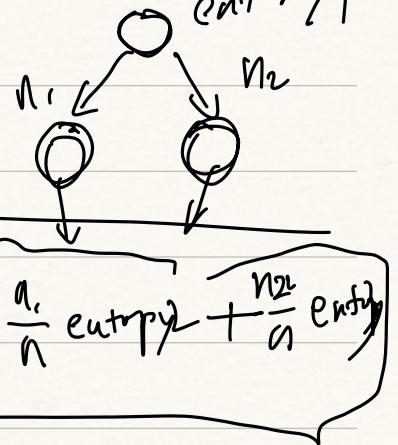


For DTs, if we have K split,

then we will have at most 2^{K-1} nodes!

maximal number in the last level

DT → universal approximator like Neural Network
(universal approximation theorem)

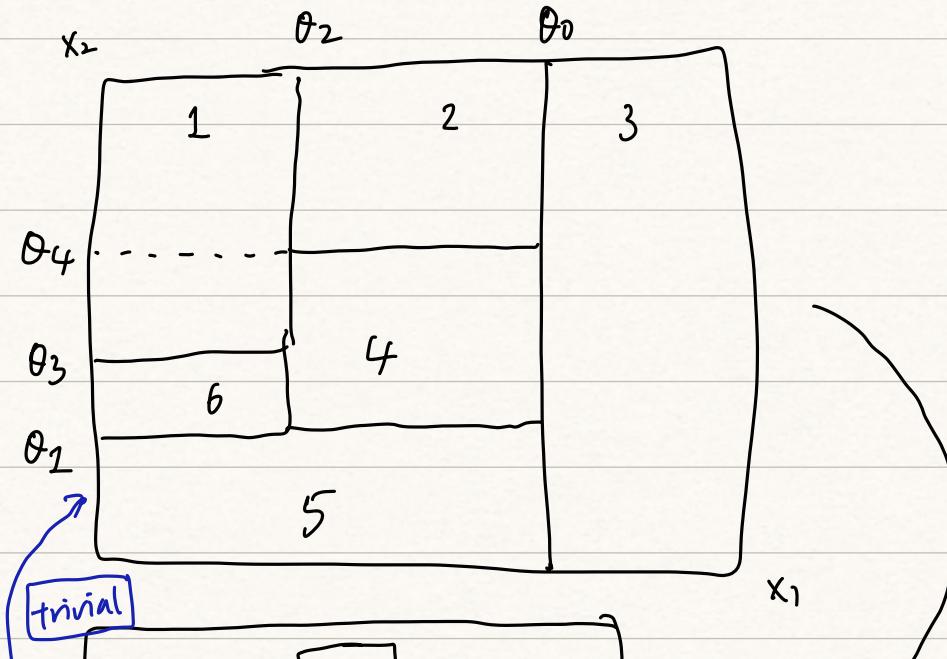


Formalize

★ a general decision tree, Hypothesis space is:

$$H = \{ f: f(x) = \sum_{j=1}^J a_j \mathbb{1}_{\{x \in R_j\}} ; \{R_j\}_{j=1}^J \text{ is a partition of } X, a_j \in Y \}$$

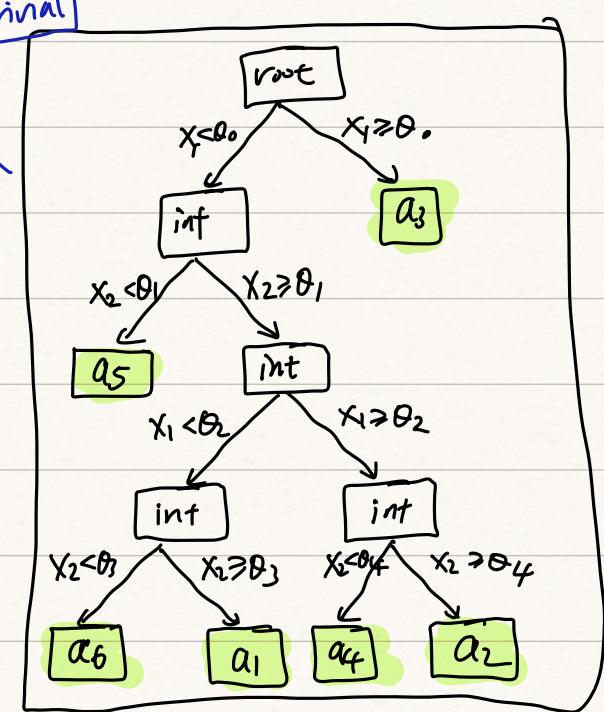
Target (Oracle) $f^*: x \in X \mapsto y \in Y$



DT Diagram

(from partition graph)

equivalent!



key trick: 找一条到底的线

Question: How to learn DT? (optimal) parameter a_j & R_j)

Ans

Optimal a_j

Easy Part

→ For **Regression Problem**, \hat{a}_j = average target over region R_j

achieved by minimizing
square loss

$$= \frac{\sum_i y_i \mathbb{1}_{\{x_i \in R_j\}}}{\sum_i \mathbb{1}_{\{x_i \in R_j\}}}$$

$$= \frac{\sum_{x_i \in R_j} y_i}{|R_j|}$$

= average target within region R_j

→ For **classification**, majority wins

Difficult Part

Optimal partition $\{R_j\}_{j=1}^J$

$\Leftrightarrow a_j = \text{mode } \{y_i : x_i \in R_j\}$

$$\frac{1}{2} \sum_{j=1}^J \sum_{x_i \in R_j} (y_i - a_j)^2$$

→ **Regression**

optimization problem

minimize empirical risk

$$\min_{\{R_j\}} \frac{1}{2} \sum_{i=1}^N (f_{\text{or}}(x_i) - y_i)^2$$

forms partition

$$= \min_{\{R_j\}} \frac{1}{2} \sum_{i=1}^N \left(\sum_{j=1}^J a_j \mathbb{1}_{\{x_i \in R_j\}} - y_i \right)^2$$

find good-enough solution

trade-off

find optimal sol \approx

computationally infeasible

(all possible partitions)

since partition number

(Bell number $B(n)$)

$$B(n) \gg e^n$$

NP-Hard

greedy solution

Recursive Binary Splitting

* always is not optimal

① pick one dimension (randomly or ...)

② find optimal split θ to split the input

dimension

$$\begin{cases} x \geq \theta \\ x < \theta \end{cases}$$

w.r.t some optimization problem

Binary split with θ

???

$$-\sum_{j=1}^J \sum_{k=1}^K p(k|j) \log p(k|j) \text{ reasonable?}$$

Classification

Modify the loss function

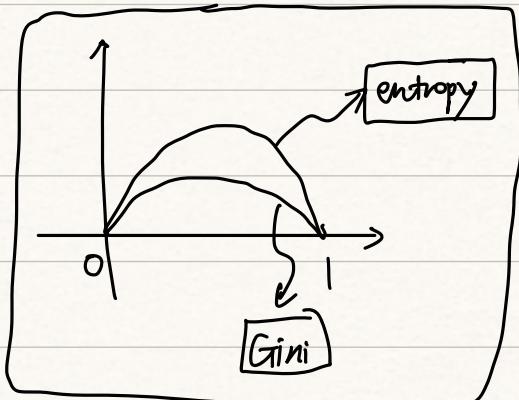
① Entropy:

$$-\sum_{j=1}^J \sum_{k=1}^K p(k|j) \cdot \log p(k|j)$$

② Gini Impurity:

$$\sum_{j=1}^J \left[\sum_{k=1}^K p(k|j) (1 - p(k|j)) \right]$$

$\rightarrow P(\text{In region } j, \text{ select 2 different classes})$



a.) $p(k|j) \rightarrow$ proportion that Region j

data belongs to class k

$$b) \sum_k p(k|j) = 1$$

c) achieves minima when

$$p(k_j|j) = 1 \text{ for } \forall j \in J$$

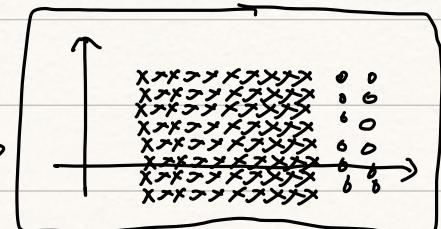
the most pure case!

Disadvantage

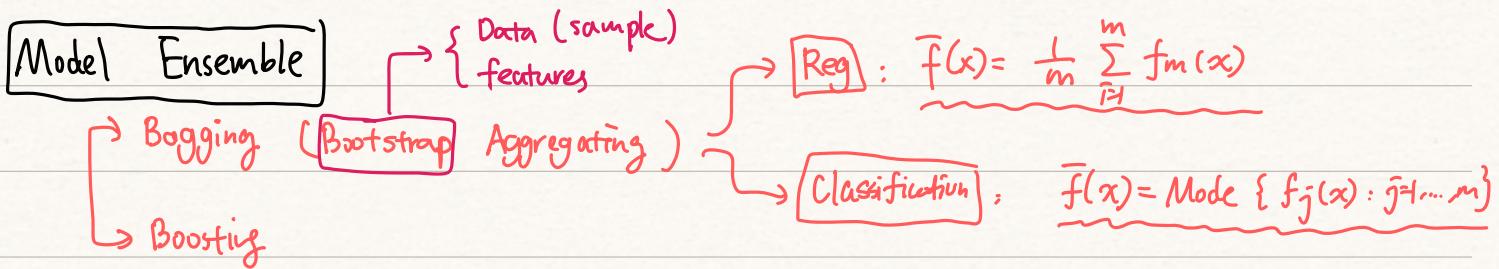
→ greedy algo \Rightarrow sub-optimal

→ sensitive to data variation & imbalanced data

→ over-fitting



Model Ensemble (Bagging)



Bagging

$$f_i(x) = f^*(x) + \underbrace{\varepsilon_i(x)}_{\text{noise}}$$

Assume

$$\left\{ \begin{array}{l} \mathbb{E} [\varepsilon_i(x)] = 0 \rightarrow \text{unrealistic} \\ \text{Var} [\varepsilon_i(x)] = \sigma^2(x) \\ \text{Cov} [\varepsilon_i(x), \varepsilon_j(x)] = 0 \end{array} \right.$$

Define Error

$$\textcircled{1} \quad E(x) = \mathbb{E} [f_j(x) - f^*(x)]^2$$

$$\frac{1}{m} \sum_{j=1}^m [f_j(x) - f^*(x)]^2 = \mathbb{E} [\varepsilon_j(x)]^2$$

$$\frac{1}{m} \sum_{j=1}^m \mathbb{E} [\varepsilon_m(x)^2] = \sigma^2(x)$$

$$\textcircled{2} \quad \bar{E}(x) = \mathbb{E} [\bar{f}(x) - f^*(x)]^2$$

$$= \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m (f_i(x) - f^*(x)) \right]^2$$

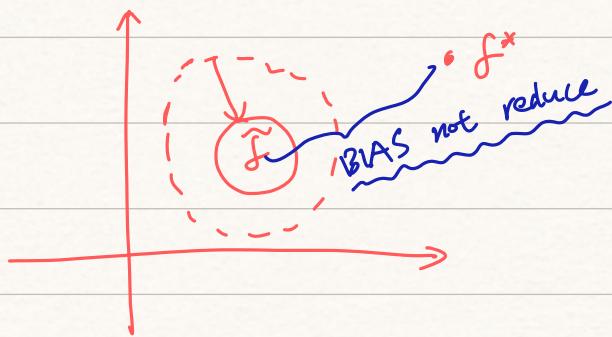
$$= \frac{1}{m^2} \mathbb{E} \left[\sum_{i=1}^m (f_i(x) - f^*(x))^2 + \sum_{i \neq j} (f_i(x) - f^*(x)) * \sum_{j \neq i} (f_j(x) - f^*(x)) \right]$$

$$= \frac{1}{m^2} \mathbb{E} \left[\sum_{i=1}^m \varepsilon_i(x)^2 + \sum_{i \neq j} \varepsilon_i * \varepsilon_j \right]$$

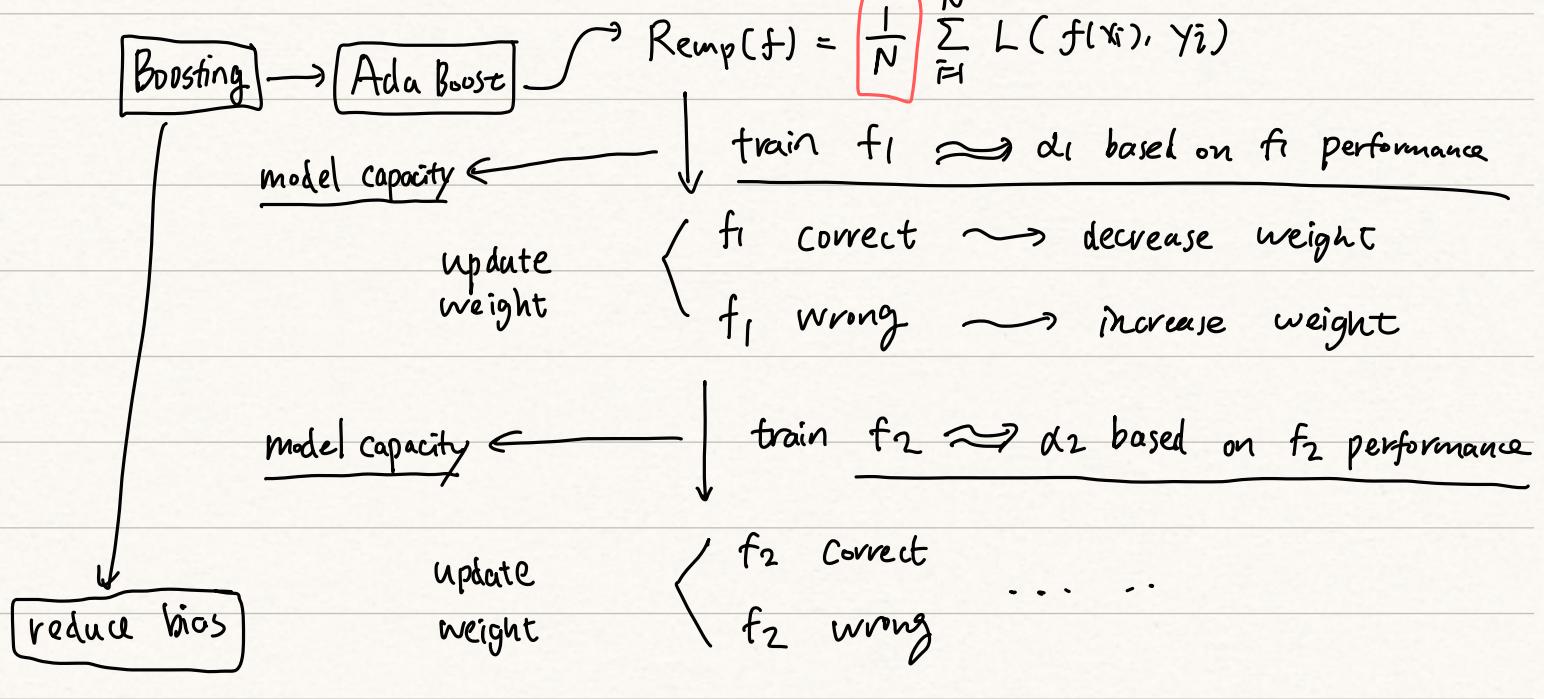
$$= \frac{1}{m} \sigma^2(x)$$

$$= \frac{1}{m} E(x) \rightarrow \boxed{\text{shrink}}$$

Bogging → reduce variance



uniform weight (unweighted)



→ Aggregate Model f_1, \dots, f_n

& model capacity $\alpha_1, \dots, \alpha_n$

Cross-Validation [CV] → choose hyper-parameters

① use validation set { train validate test

② cannot use test data to select

model

Solution (CV)

**BAGGING
IDEA**

CROSS-validation

a) split it into training & test

b) split training set into K folds

c) train model for K times

(i -th validation, rest training) $i=1 \dots K$

\downarrow
 $score_i$ $i=1, 2, \dots, K$ error on validation set

d) average score $_i$ $i=1, 2, \dots, K$

\downarrow performance of one

specific choice of hyper-param

\downarrow
Score (hyper-param)

e) compare score over different

choices of hyper-parameters

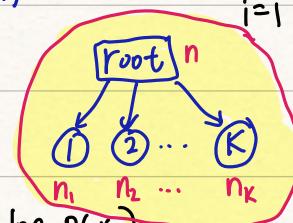


determine the optimal hyper-parameter!

My Personal Perspective

→ The Idea of GAIN - SPLIT := Entropy (root) - $\sum_{i=1}^k \frac{n_i}{n} \text{Entropy}(i)$

\downarrow
if split to K nodes



Recall: ① entropy $H(X) = - \sum p(x) \log p(x)$

② joint entropy $H(X, Y) = - \sum_{x,y} p(x, y) \log p(x, y)$

$H(X, Y) \geq H(X)$

③ conditional entropy : $H(X|Y) = \sum_y P(y) H(X|Y=y)$

$E_y [H(X|Y=y)] = - \sum_y p(xy) \log p(x|y)$

A more reasonable objective

for classification DT:

[objective]

$$\min_{\text{partition}} H(X | \text{Partition})$$

$$\Leftrightarrow \min_{\substack{\{D_1, \dots, D_J\} \\ \{\text{partition}\}}} \sum_{j=1}^J \sum_{k=1}^K p(c_k, j) \log p(c_k | j)$$

then, it is reasonable to

consider [gain-split] for each split

$$\text{GAIN_SPLIT} = \text{Entropy}(\text{root}) - \sum_{i=1}^T \frac{n_i}{n} \text{Entropy}(i)$$

Normalization

$$\text{Gain_split_ratio} = \frac{\text{Entropy}(\text{root}) - \text{Entropy}(\text{root}|\text{split})}{\text{Entropy}(\text{split})}$$