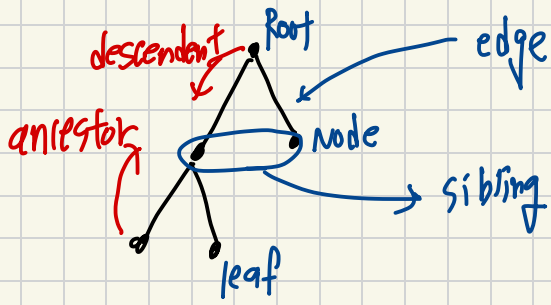


# < Tree >



## < Regression Tree >

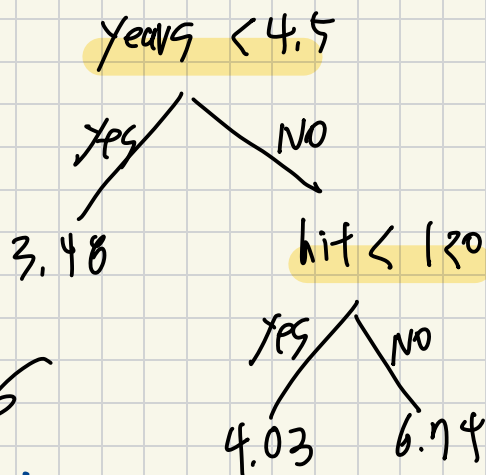
$X \in \mathbb{R}^d$  인 vector 이다.

우선 매우 간단하게  $X \in \{Hit, Year\}$  2개의 variable을 갖는 경우 고려.

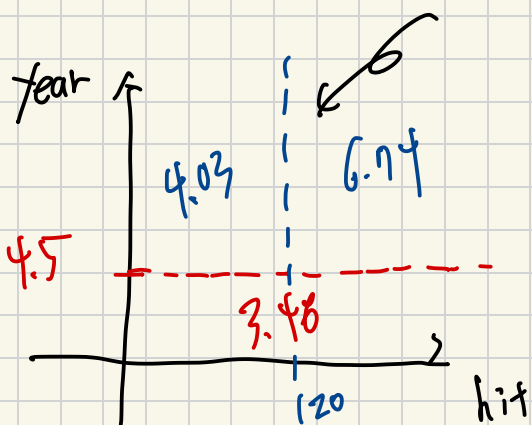
예를 들어  $|Hit| = 9140, 20$ 의 범위를 표현이 가능하다.

그러므로 tree의 node의 경우  $x_j < t_k$  등으로 표현이 된다.

example)



internal node : 구간분  
leaves node : 결과.



leaves node 수 많을수록 split 된다.

→ 어느정도의 해석가능성을 제공한다.

## < Building Regression Tree >

최종 목표: Non-overlapping 하게 Area를  $R_1, R_2 \dots R_N$  으로 나눈다.

같은 영역 값들은 모두 Same 상수로 결정.  $\rightarrow$  너무 Simple?

만일 더 detail한 Regression 필요시 leaves node를 늘리면 된다.  $\rightarrow$  해결.

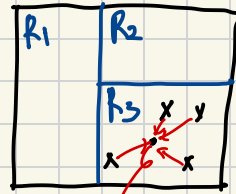
Regression을 주는 법

$$\min \sum_{j \in R} \sum (y_j - \hat{y}_j)^2 \text{ 을 만족해야 한다.}$$

그래서 각 Region마다 그 차이를 minimize 하는 대표값  $C$  설정.

$$\hat{y}_j = \arg \min_c \sum_{i \in R_j} (y_i - c)^2$$

$\hookrightarrow$  ML에 준하면 결국 R이 포함된 data들의 평균이 나오게 된다.



가장 잘 represent 하는  $C$

평균은 데이터들의 오차 제곱을 최소화 해주는 하나의 대표값이다.

절대값의 오차를 최소화 해주는 대표값은 중앙값이다

## < Top-down Greedy Algorithm >

노드에서 부터 재귀적 최선의 selection을 반복해 가는 방법으로 tree를 생성한다.

그러면 어느 variable  $X$ 를 어느 기준  $\sigma$ 로 비교해서 나눌지를 정해야 한다.

$\Rightarrow$  모든  $X \in \{X_1, X_2 \dots X_p\}$ 에 대해 모든  $\sigma$ 를 고려해서

$$\sum_{i \in R_1} (y - \hat{y}_{R_1})^2 + \sum_{i \in R_2} (y - \hat{y}_{R_2})^2 \text{ 을 가장 작게 하는 } X_1 \text{ 선택.}$$

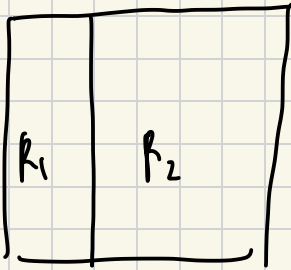
4 설명 (1)  $X_j$ 를 오름차순 정렬

(2)  $X_{jt}$ 과  $X_{j,t+1}$  사이의 공간값을  $s_t$ 로 사용

(3) 무등호는  $\leq s$  or  $< s$  사용 (분리하기만 하면 된다)

\* 주로 1개의 Region에  $K$ 개의 data가 있을 때 가지 뻗어간다.

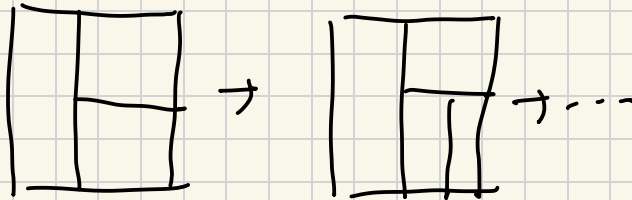
(1)



(2)  $R_1$ 을 split한거  $R_2$ 을 split 한것도

평평  $\Rightarrow$  모든 평평을 다 따져서 greedy 하기

$n \rightarrow n+1 \rightarrow \dots$  이기 Region을 늘려나간다.



## < Tree pruning >

Region 위에 data point가  $\downarrow$  이면 overfitting이 발생한다.

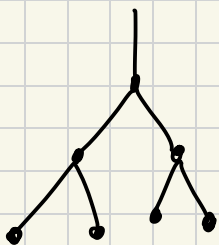
$\rightarrow$  Boundary를 smooth하게 (작은 번만 split)  $\rightarrow$  Tree를 greedy로 많이 만들어

Model이 성능이 좋아질지 모름 (훈련 다 해봐도 model 정확도 파악 가능)

그래서 우선 split 한 후 & 다시 Region을 늘려가는 방법이 pruning.

$$\sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$

$\rightarrow$  leaves node가 많을수록

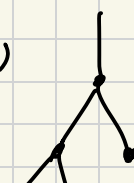


$\Rightarrow$

(1)

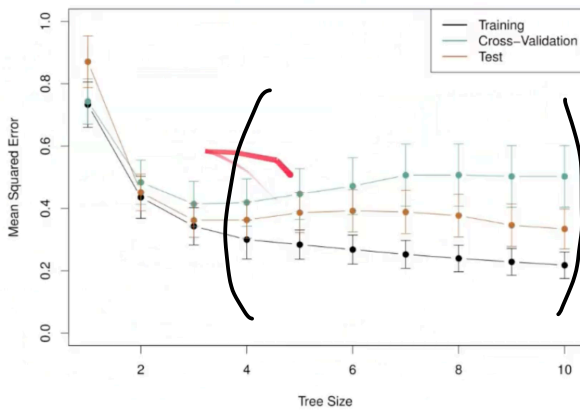


(2)



(1), (2)를 compare.

각 대표값 1개 갖아서 (2 group 정도)  
그 중  $R$ 가 min인 features  
node를 merge 한다.



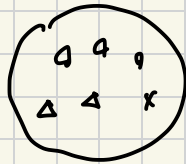
→ over fitting 자증.

즉, leaves node = 3이 될 때까지 pruning 진행.

## < Classification Tree >

부가 (categorical) ⇒ Region은 나무를 select하는 방법.

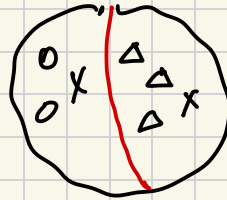
- (1) Region split하기  
 (2) Representation value하기



o:  $\frac{8}{8}$  Δ:  $\frac{2}{8}$  x:  $\frac{1}{8}$  을 위해서 o가 하도.

Split시 misclassification ratio를 minimize.

"Impurity" : Region에 포함된 symbol class가 많지 않도록.



minimize)

$$\frac{1}{N_1} \text{Impurity}(R_1(V, S)) + \frac{1}{N_2} \text{Impurity}(R_2(V, S))$$

## < Impurity >

$$\text{Impurity}(R) = - \sum_{k=1}^n \underline{p_{j,k}} \log p_{j,k}$$

$$\text{Impurity}(R_j) = - \sum_{k=1}^n p_{j,k} \log p_{j,k}$$

○ This is called **entropy**.

○ Illustration by example:

- [1, 0, 0, 0, 0] → -1 log 1 = 0 // 만일
- [0.5, 0.5, 0, 0, 0] → -2 \* 0.5 log 0.5 = log 2
- [0.2, 0.2, 0.2, 0.2, 0.2] → -5 \* 0.2 log 0.2 = log 5

↓ 불순도 증가.

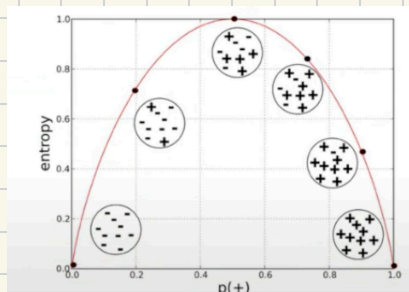
## < Gini-index >

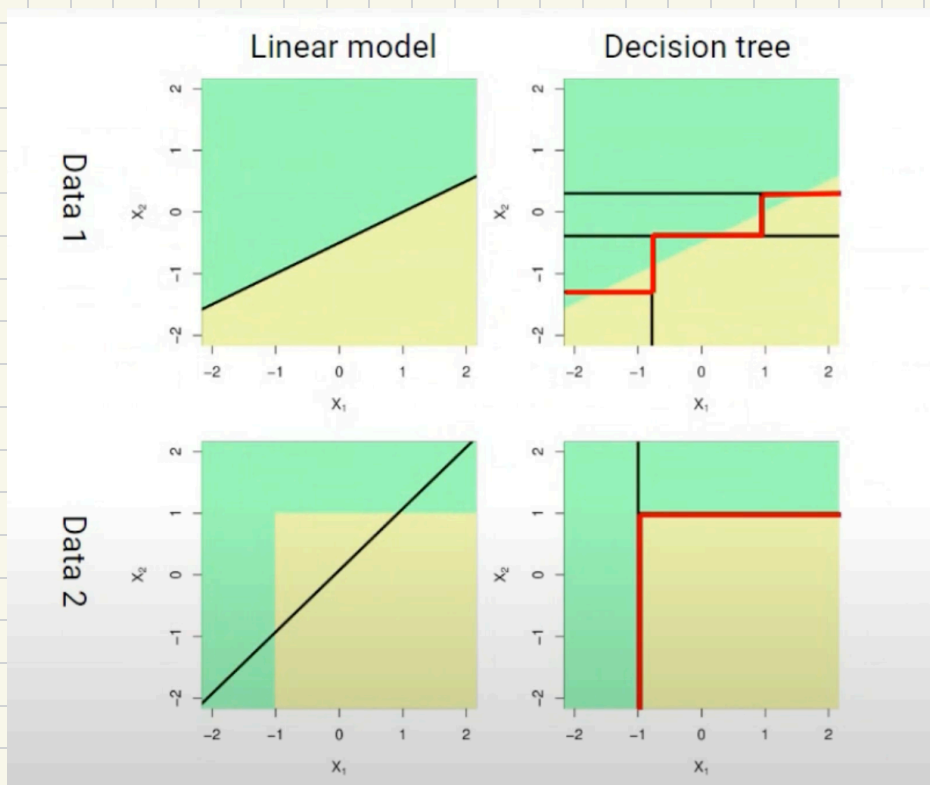
$$G = \sum_{j=1}^K p_{j,k}(1 - p_{j,k})$$

○ Illustration by example:

- [1, 0, 0, 0, 0] → 0
- [0.5, 0.5, 0, 0, 0] → 2 \* 1/4 = 0.5
- [0.2, 0.2, 0.2, 0.2, 0.2] → 5 \* 0.2 \* 0.8 = 0.8

○ Gini index and the entropy are very similar numerically.





심지어는 Data가 원래  
분류와 비슷하게 더 잘된다.

But Decision Tree가 더  
비선형적 특징을 배린다

But Tree는 axis align style의  
관점이 존재한다.

(강한점이 존재)

Tree의 pro : (1) 설명하기 쉽고  
(2) 시각화가 가능하리.

Tree의 cons : (1) 성능이 좋지 않다.

(2) Robust 하지 않다  $\rightarrow$  Train에만 잘 작동....! (치명적)

$\rightarrow$  Data가 추가되면 tree의 모양이 너무 많이 change 된다.