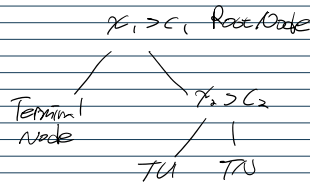


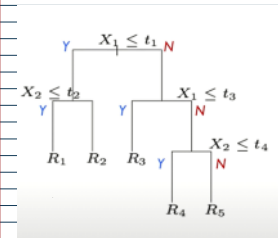
# 1. Decision Tree

데이터가 어떤 범주에 속하는지 결정하기 위한 방법.

- 분기: 새로운 변수를 기준으로 분기
- 예측: 분기된 후의 결과



# 2. Regression Tree



$C_m$ : Regression Tree에서 예측된  $R_m$ 의 평균값 (mean...)

$$\hat{f}(x) = \sum C_m I\{(x_1, x_2) \in R_m\}$$

$$I: \text{Indicator func.} \begin{cases} 0 \rightarrow F \\ 1 \rightarrow T \end{cases}$$

① 데이터  $M$ 에 대한 분할  $f(x) = \sum C_m I(x \in R_m)$

② 회귀분석은 cost function을 최소화.

$$\min_{C_m} \sum (y_j - f(x_j))^2 \quad C_m = \text{ave}(y_j | x_j \in R_m)$$

→ 각각의  $y$ 값 평균을 예측하기 위한 것.

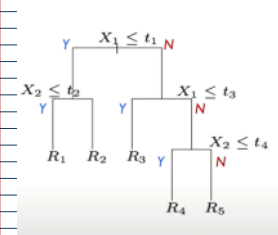
분할 기준을  $j$ 로 정함

$$R_1(j, s) = \{x | x_j \leq s\} \quad R_2(j, s) = \{x | x_j > s\}$$

$$\arg \min_{j, s} \left[ \min_{C_1} \sum_{x_j \in R_1} (y_j - C_1)^2 + \min_{C_2} \sum_{x_j \in R_2} (y_j - C_2)^2 \right]$$

$$\sum_{x_j \in R_1} (y_j - C_1)^2 + \sum_{x_j \in R_2} (y_j - C_2)^2 \quad \leadsto j, s \text{을 찾아야 하는 것}$$

# 3. Classification



→ 정답이 맞는지 틀린지 판별

$\hat{p}_{mk}$ : 클래스  $m$ 에서  $k$ 클래스에 속한 관측치의 비율

$$= \frac{1}{N_m} \sum I(y_j = k)$$

클래스  $m$ 의 분류된 관측치는  $k(m)$ 클래스로 분류.  $k(m) = \arg \max_k \hat{p}_{mk}$

$$\hat{f}(x) = \sum k(m) I\{(x_1, x_2) \in R_m\}$$

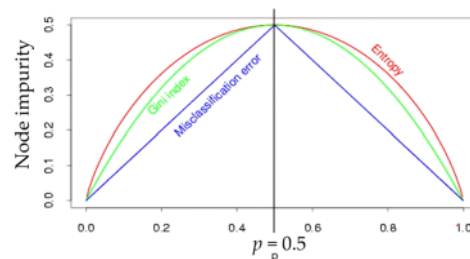
분류 성능: 측정

① Misclassification Rate  $\left[ \frac{1}{N_m} \sum I(y_j \neq k(m)) = 1 - \hat{p}_{mk(m)} \right]$

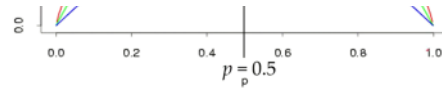
→ 실제 정답과 예측한 답이 다를 때.

② Gini Index  $\left[ \sum_{k \neq k'} \hat{p}_{mk} \hat{p}_{mk'} = \sum_{k=1} \hat{p}_{mk} (1 - \hat{p}_{mk}) \right]$

③ Cross-Entropy  $\left[ - \sum_{k=1} \hat{p}_{mk} \log \hat{p}_{mk} \right]$



$$\textcircled{3} \text{Cross-Entropy} \left[ -\sum_{k=1} p_{mk} \log p_{nk} \right]$$



$\Rightarrow$  노드에서 분할 및 분할 후의 분할 값의 차이 정도

Information Gain : IG : 분할 후 Entropy 감소.

$$IG(S, A) = Entropy(S) - \sum_{v \in V(A)} \frac{|S_v|}{|S|} Entropy(S_v) \quad : \text{원 A에 의한 IG.}$$

#### 4. 정보의 효율 측정

- 계층적 : 층이 여러 보일수록 다음 단계로 전파.
- 노드가 많음
- 과소적합