

- Decision Tree.

- classification and Regression

둘 다 가능!

- 예측 결과물에 의해 쉽게 쉽게 설명할 수 있는 자주 사용하는 모델이다.

(The output is a set of rules)

- 설명 변수의 형태 상관없이 사용 가능.

(Able to handle both numeric, cat)

- Decision Tree (classification)

- CART Algorithm

→ Classification & Regression Tree

\approx CHAID, C4.5, C5.0

- Recursive Partitioning

→ 순차적으로 분할한다.

(Impurity를 기준으로)

- Pruning the Tree

→ Full Tree or Merge (Post-Pruning)

사전 조건으로 Cut. (Pre-Pruning)

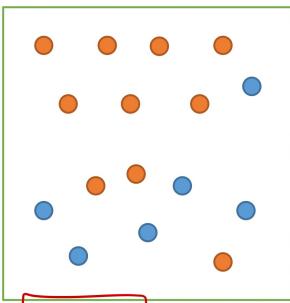
• Recursive Partitioning

- Impurity 를 기준으로 분할 한다.

① Gain Index.

→ CART에서 사용하는 measure

$$\rightarrow I(A) = 1 - \sum_{k=1}^m p_k^2$$



$$\begin{aligned}
 I(A) &= 1 - \sum_{k=1}^m p_k^2 \\
 &= 1 - \left(\frac{6}{16}\right)^2 - \left(\frac{10}{16}\right)^2 \\
 &\approx 0.47 \quad \Rightarrow \text{Impurity가 높은 편.}
 \end{aligned}$$

• 개수 : 10 , ● 개수 : 6 .

$$\rightarrow \text{Max} : 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = \frac{1}{2} = \underline{\underline{0.5}}$$

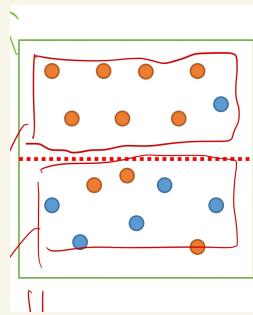
$$\rightarrow \text{Min} : 1 - 1^2 - 0^2 = \underline{\underline{0}}$$

→ 즉, 작을수록 좋은 분할 .

Gini Index를 사용하여 가장 좋은 분할
방법 찾기.

$$\rightarrow I(A) = \sum_{i=1}^d \left(R_i \left(1 - \sum_{k=1}^m P_{ik}^2 \right) \right)$$

\Rightarrow Gini Index의 가중합.



$$I(A) \quad \hookrightarrow \text{설정된 데이터의 비율.} \quad \downarrow \\ = 0.5 \times \left(1 - \left(\frac{7}{8} \right)^2 - \left(\frac{1}{8} \right)^2 \right) + 0.5 \times \left(1 - \left(\frac{3}{8} \right)^2 - \left(\frac{5}{8} \right)^2 \right) \\ = 0.34$$

\rightarrow Information Gain

$$= 0.49 - 0.34 = \underline{\underline{0.14}}$$

\Rightarrow split에 대한 효과도 해석 가능.

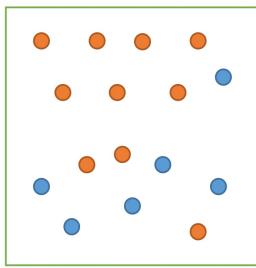
\Rightarrow split point에 대해 Information Gain

을 비교하여 가장 큰 값을 가진
split point 를 찾는다!

② Deviance.

→ 대푯값 Tree의 사용하는 것.

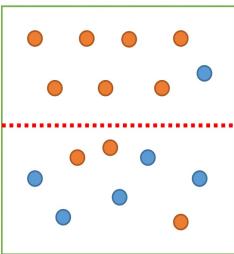
$$D_i = -2 \sum_{k=1}^m n_{ik} \log(P_{ik})$$



$$D_i = -2 \times \left(10 \times \log\left(\frac{10}{16}\right) + 6 \times \log\left(\frac{6}{16}\right) \right)$$

$$= 21.17$$

→ Min : 0 (same class)



$$D_1 = -2 \times \left(7 \times \log\left(\frac{7}{8}\right) + 1 \times \log\left(\frac{1}{8}\right) \right) = 6.03$$

$$D_2 = -2 \times \left(3 \times \log\left(\frac{3}{8}\right) + 5 \times \log\left(\frac{5}{8}\right) \right) = 10.59$$

$$D_1 + D_2 = 16.62$$

Information gain

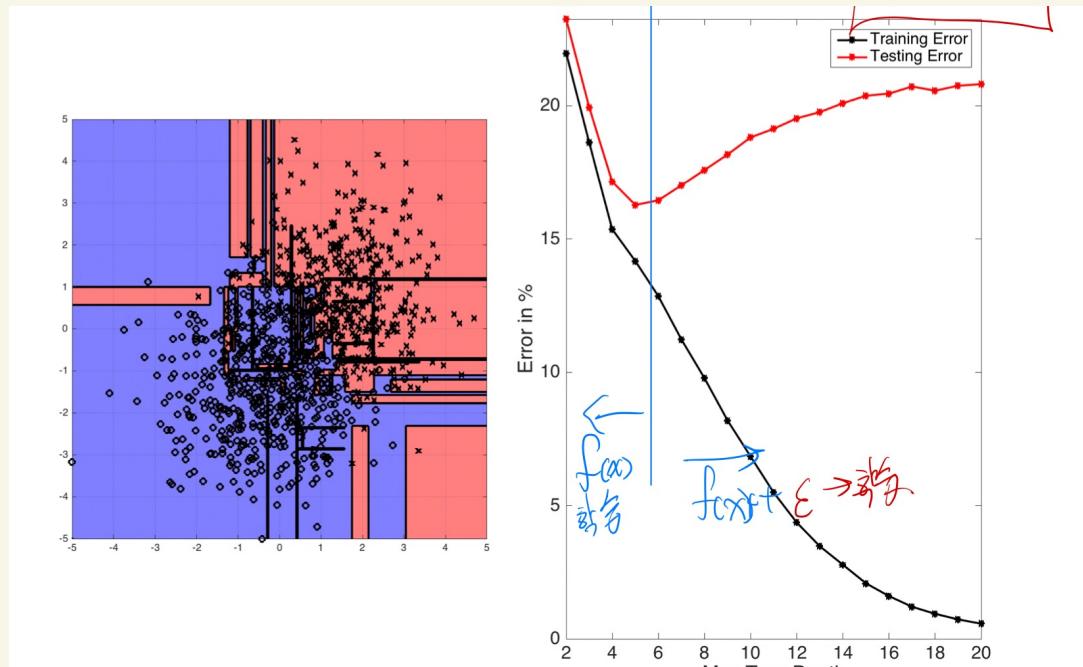
$$= 21.19 - 16.62 = \underline{\underline{4.55}}$$

\Rightarrow split ~~not~~

• Full Tree 만족

- 설명변수 : m .
- recode : n .
- 불활 순간마다. m 개의 변수에서
Sorting 후 불활 포인트를 찾고,
Information gain을 구한 다음
최대의 값을 가지는 split을
한다.
- 중요조건 : Impurity 가 0이 되는
순간 leaf Node로 설정 후
더 이상 불활 X.

- Full Tree \Rightarrow 100% of Acc \Rightarrow 100% overfitting
 이 모델이 좋지 않다!



· Pruning.

→ Overfitting 문제 / 과적합
의 해:

→ Merge.

→ Cost complexity 가로드!

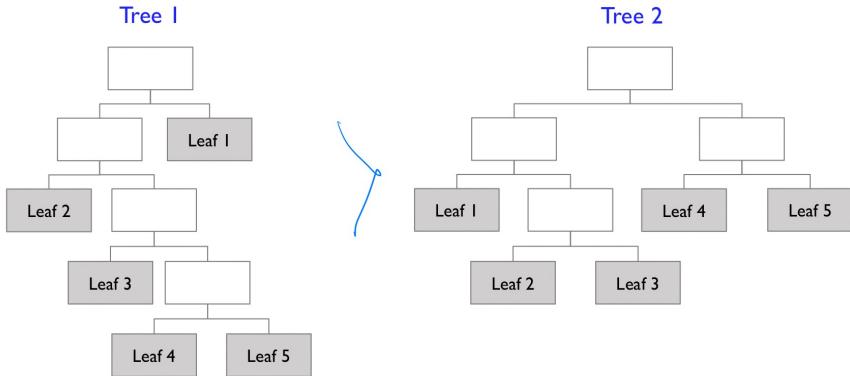
$$CC(T) = Err(T) + \alpha \times L(T)$$

* $Err(T)$

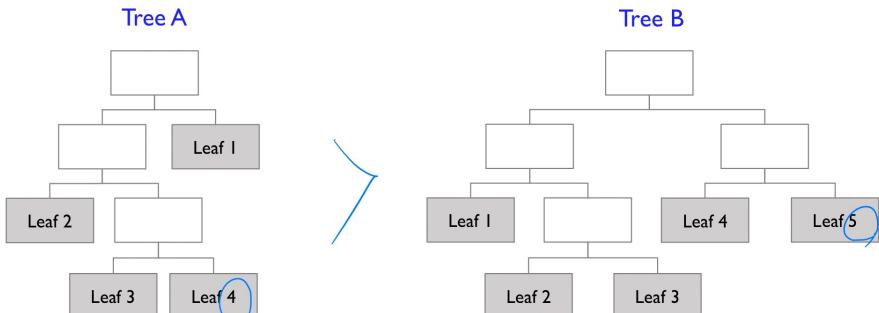
→ proportion of misclassified
records in the variables data

$\#L(T)$: Number of leaf Node

Cost Complexity Example 1

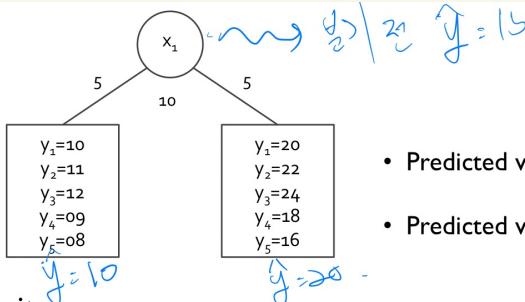


Cost Complexity Example 2



• Regression Tree

→ output: 정답의 예측값!



- Predicted value of the left leaf node = 10
- Predicted value of the right leaf node = 20

• Impurity

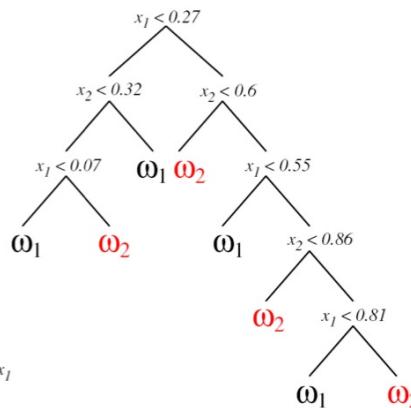
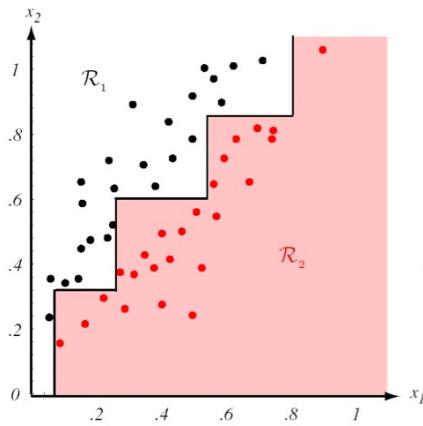
✓ Sum of squared error (SSE: $\sum_{i=1}^n (y_i - \hat{y})^2$)

✓ SSE(Parent) = 300, SSE(Left) = 10, SSE(Right) = 40, Gain = 250

분산 → 분산

Information Gain //

CART Disadvantages.



One variable at a time

