

Decision Tree

중간 범위X

Regression
이전에 배운 Logistic Regression
Naïve Bayesian
과는 달리 수학적 개념이 사용X

Tennis Playing

• Tennis playing record

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Tennis Playing

- ➊ Today is Sunny, Mild, High and Strong
 - Will he play tennis today?
- ➋ How can you predict his tennis playing?
- ➌ What is the rule for playing tennis?
 - How to extract patterns from the data?

Basic Idea

→ Do not consider all inputs
1. Let's choose ~~Just one important input~~ Just one important input Decision
2. predict based on the chosen input

Input이 (7U의) 어떤 능력

- We have a simple data

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

Input

APL이면 $A(T) \rightarrow \text{Play}(T)$
 $A(F) \rightarrow \text{play}(F)$

- Will he play tennis if $A=F, B=T, C=T$?
- Will he play tennis if $A=T, B=T, C=F$?

① choose
②

Basic Idea

- A simple idea: See one of inputs and answer based on the majority

How evaluate to How Important?

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

A=T



Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T

2T OF
100%

A=F



Day	A	B	C	Play
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

1T 2F
67%
not bad

Basic Idea

- A simple idea: See one of inputs and answer based on the majority

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

$B=T$

Day	A	B	C	Play
3	F	T	F	T

100%

$B=F$

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
4	F	F	T	F
5	F	F	F	F

2T2F
50%

Not useful
A>T C<Cn!

Basic Idea

- A simple idea: See one of inputs and answer based on the majority

(Good, Not Bad) (Good, Bad) Bad, Not Bad

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

$$C=T$$



Day	A	B	C	Play
1	T	F	T	T
4	F	F	T	F

Impure set

T/F 50% 50%

1 T 1 F
50%

$$C=F$$



Day	A	B	C	Play
2	T	F	F	T
3	F	T	F	T
5	F	F	F	F

2 T 1 F

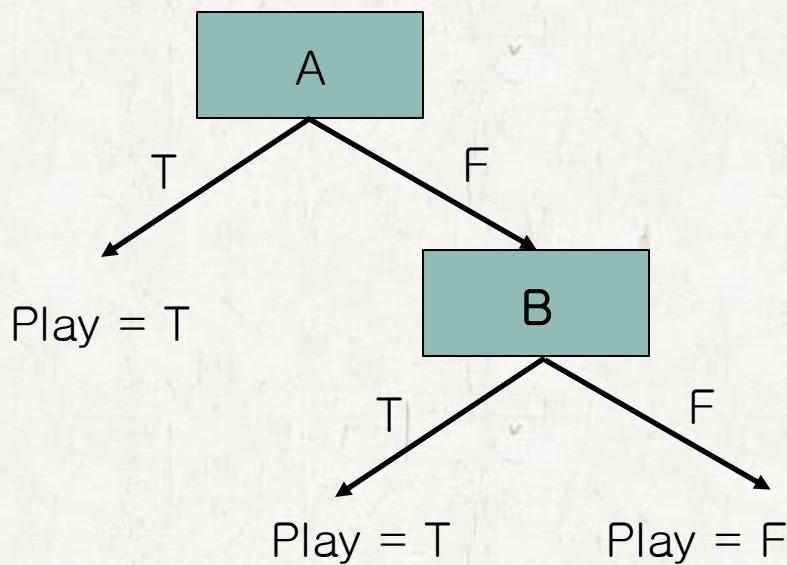
67%

Not Bad

Basic Idea

- A simple idea: See one of inputs and answer based on the majority

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F



- ① measure Impurity of each subset
- ② measure Impurity of split
- ③ Evaluate Gain of split
→ 풀수록 좋다
- ④ choose maximum Gain

choosing a Input

Attribute Selection for Split

Impurity measure

Which Attribute Is the Best?

Day	A	B	C	Play
3	F	T	F	T

B=T

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

B=F

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T

A=T

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T

$F \rightarrow T$ pure set
 $T \rightarrow T$

Day	A	B	C	Play
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

$F \rightarrow T$ less impure
 $F \rightarrow T$ impure
 $F \rightarrow F$

Day	A	B	C	Play
1	T	F	T	T
4	F	F	T	F

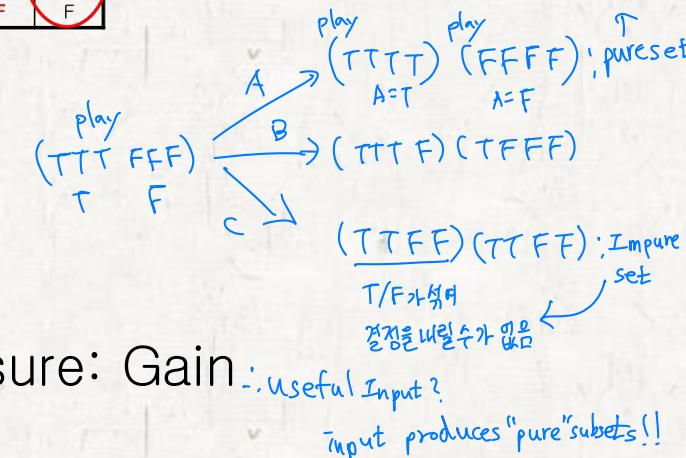
C=T

Day	A	B	C	Play
2	T	F	F	T
3	F	T	F	T
5	F	F	F	F

A=F

decision 가능!
(T나 F밖에 없음)

- We will define
 - Impurity of single sets: Entropy
 - Impurity of splits: Average of Entropy
- Then, present attribute selection measure: Gain

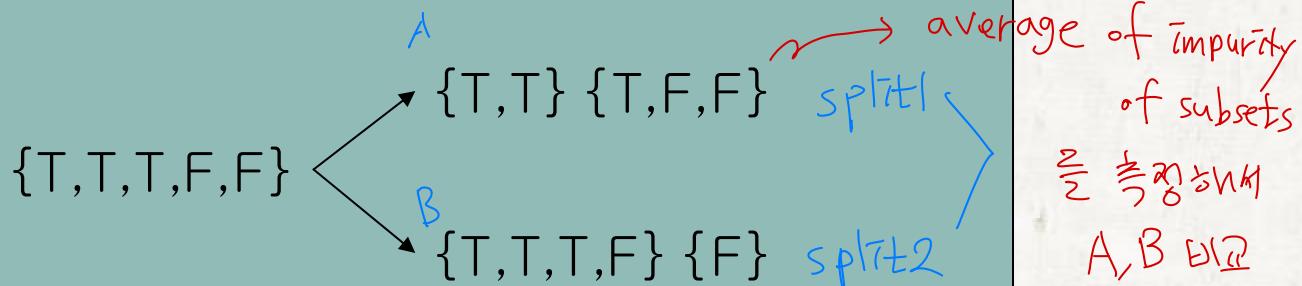


Impurity Measure

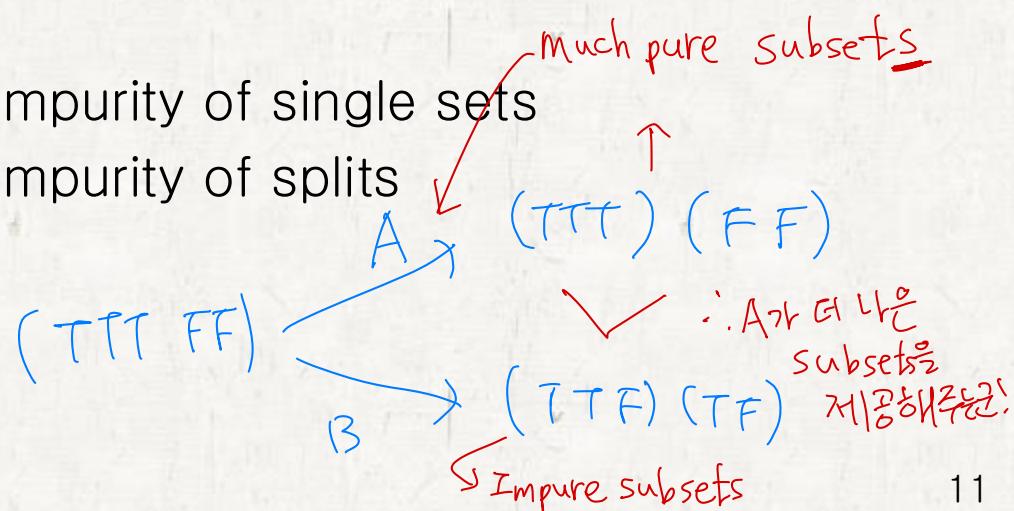
- We need the impurity of splits

데이터가 많아지면
Impurity 측정 힘들
→ 측정방법을 고민하자

Which split is composed of purer subsets?



- Step1: We will define the impurity of single sets
- Step2: We will define the impurity of splits



4step 중에 ①

↳ Impurity Measure for Single Sets

- Let's define impurity measure for single sets

$F(\text{Set}) = \text{Impurity of Set in } [0,1] \text{ (portion of element)}$

구질성에 따지는 함수

Set₁: (1, 1, 1, 1, 1, 1, 1, 1)

Set₂: 1, 1, 1, 1, 1, 1, 1, 0

Set₃: 1, 1, 1, 1, 1, 1, 0, 0

Set₄: 1, 1, 1, 1, 1, 0, 0, 0

Set₅: 1, 1, 1, 1, 0, 0, 0, 0

Set₆: 1, 1, 1, 0, 0, 0, 0, 0

Set₇: 1, 1, 0, 0, 0, 0, 0, 0

Set₈: 1, 0, 0, 0, 0, 0, 0, 0

Set₉: 0, 0, 0, 0, 0, 0, 0, 0

~~ure for single sets~~ \rightarrow single element \rightarrow pure
(하나의 원소만 있는 것)

pure set (purest)

↳ 속서는 중요하지X 몇개의 데이터가 있는지?

$F(Set_1) = 0$ All elements are 1's

→ Impure set (Impurest)

$F(Set_5) = 1$ Half and half.
 Totally, mixed up

$F(Set_9) = 0$ All elements are 0's

→ pure set (purest)

\Rightarrow 함수를 만드어보자
(찾아보자) 12

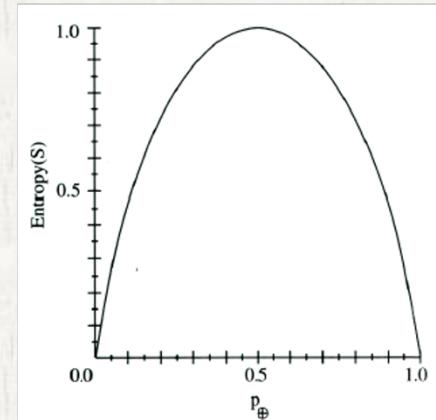
Impurity Measure for Single Sets

Impurity of single sets

$$\text{Entropy}(\text{set}) = -P_1 \log_2 P_1 - P_0 \log_2 P_0$$

P_1 = the probability that 1 appears in set

P_0 = the probability that 0 appears in set



$$A = \{0,0,0,0,0,0,0,0,0\} \quad \begin{array}{l} P_1 = 0 \\ P_0 = 1 \end{array} \quad \text{Entropy}(A) = -\frac{0}{8} \cdot \log_2 \frac{0}{8} - \frac{8}{8} \cdot \log_2 \frac{8}{8} = 0$$

$$A = \{1,1,0,0,0,0,0,0,0\} \quad \text{Entropy}(A) = -\frac{2}{8} \cdot \log_2 \frac{2}{8} - \frac{6}{8} \cdot \log_2 \frac{6}{8} = 0.81$$

$$A = \{1,1,1,1,0,0,0,0,0\} \quad \begin{array}{l} P_1 = \frac{1}{2}, P_0 = \frac{1}{2} \\ \end{array} \quad \text{Entropy}(A) = -\frac{4}{8} \cdot \log_2 \frac{4}{8} - \frac{4}{8} \cdot \log_2 \frac{4}{8} = 1$$

$$A = \{1,1,1,1,1,1,1,1,1\} \quad \text{Entropy}(A) = -\frac{8}{8} \cdot \log_2 \frac{8}{8} - \frac{0}{8} \cdot \log_2 \frac{0}{8} = 0$$

Impurity Measure for Single Sets

Impurity of single sets

- If set is composed of c_1, c_2, \dots, c_n

$$\rightarrow -P_0 \log P_0 - P_1 \log P_1 - P_2 \log P_2$$

$$Entropy(set) = -\sum_{i=1}^n P_{c_i} \log_2 P_{c_i}$$

P_{c_i} = the probability that c_i appears in set

서로 다른 요소가
모두 같은 개수로 존재할 때

Maximum

가장 많거나 적거나

more than two element

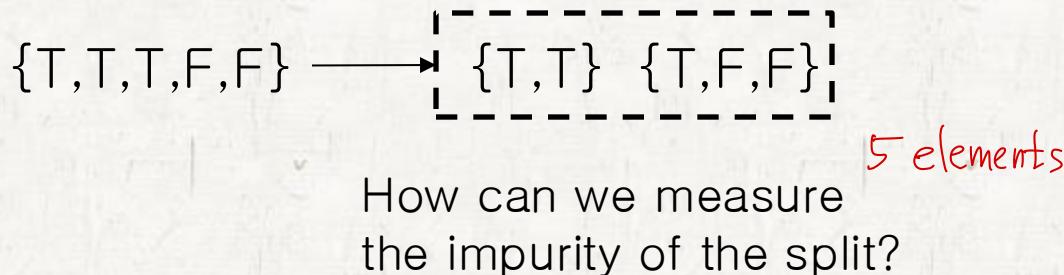
$$A = \{0, 0, 0, 1, 1, 1, 3, 3, 3\} \quad Entropy(A) = -\underbrace{\frac{3}{9} \cdot \log_2 \frac{3}{9} - \frac{3}{9} \cdot \log_2 \frac{3}{9} - \frac{3}{9} \cdot \log_2 \frac{3}{9}}$$

$$A = \{0, 0, 1, 1, 1, 1, 1, 3, 3\} \quad Entropy(A) = -\frac{2}{9} \cdot \log_2 \frac{2}{9} - \frac{5}{9} \cdot \log_2 \frac{5}{9} - \frac{2}{9} \cdot \log_2 \frac{2}{9}$$

4st step 3rd ②

Impurity Measure for Splits

• Impurity of splits



Impurity(Split) = Average of impurities of its subsets

$$\begin{aligned} \text{Impurity}(TT, TFF) &= \frac{2}{5} \times \underbrace{\text{Impurity}(TT)}_{\text{subset 1}} + \frac{3}{5} \times \text{Impurity}(TFF) \quad \text{subset 2} \\ &= \frac{2}{5} \times \left(-\frac{2}{2} \cdot \log_2 \frac{2}{2} - \frac{0}{2} \cdot \log_2 \frac{0}{2} \right) + \frac{3}{5} \\ &\quad \times \left(-\frac{1}{2} \cdot \log_2 \frac{1}{2} - \frac{2}{2} \cdot \log_2 \frac{2}{2} \right) \end{aligned}$$

Entropy

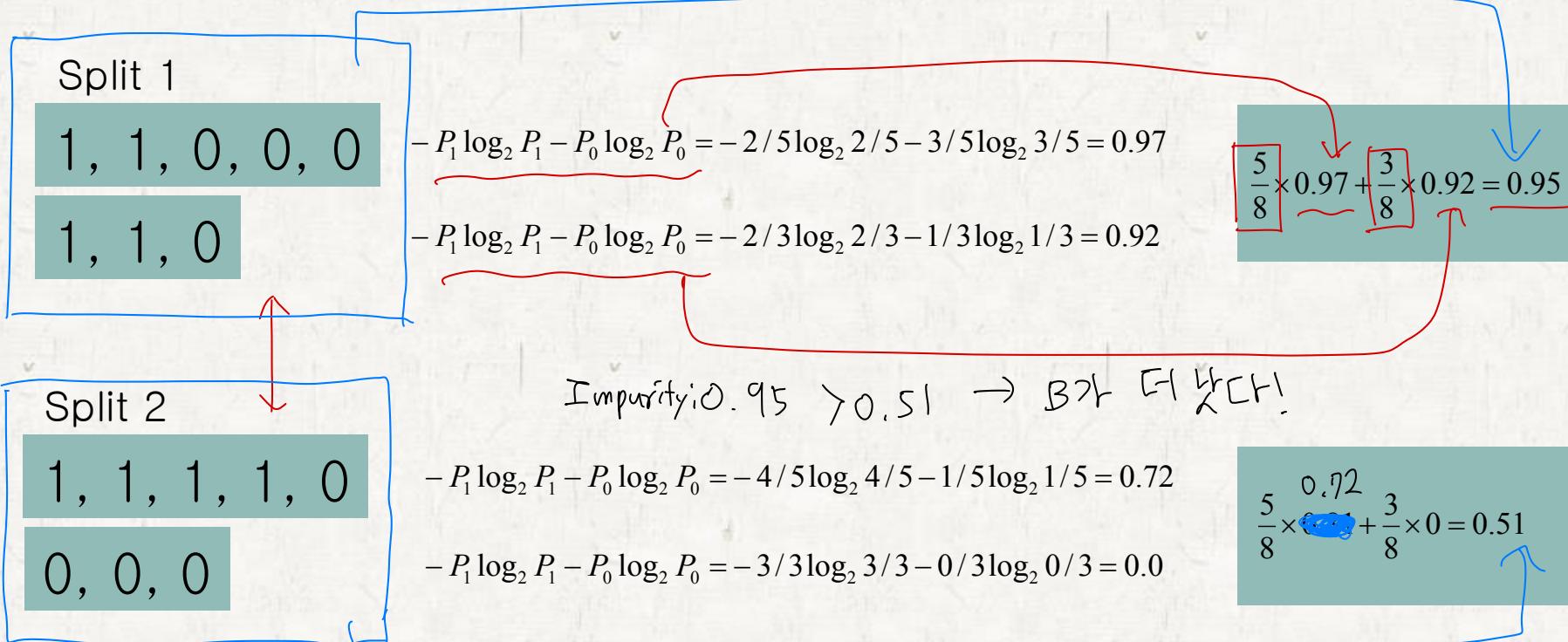
Impurity Measure for Splits

Entropy of Splits

- Example:

original set

1, 1, 1, 1, 0, 0, 0, 0



Attribute Selection

- Let's return back to a simple dataset

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

Play:{3T, 2F}

- Which attribute is good for “Split”

이전 슬라이드 참고해서 비교 후 제일 Impurity 낮은 것 choose

A=T

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

A=F

Day	A	B	C	Play
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

B=T

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T

B=F

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

C=T

Day	A	B	C	Play
1	T	F	T	T
4	F	F	T	F

C=F

Day	A	B	C	Play
2	T	F	F	T
3	F	T	F	T
5	F	F	F	F

Selecting Attribute

- Gain: Measure how much an attribute reduce the entropy

original table

$$Gain(S, a) = Entropy(S) - Entropy \text{ of splits by } a$$

chosen input to split

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

A=T

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T

A=F

Day	A	B	C	Play
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

B=T

Day	A	B	C	Play
3	F	T	F	T

B=F

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
4	F	F	T	F
5	F	F	F	F

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

C=T

Day	A	B	C	Play
1	T	F	T	T
4	F	F	T	F

C=F

Day	A	B	C	Play
2	T	F	F	T
3	F	T	F	T
5	F	F	F	F

$$Gain(Play, A) = E(TTTFF) - E(TT, TFF)$$

Entropy

$$Gain(Play, B) = E(TTTFF) - E(T, TTFF)$$

$$Gain(Play, C) = E(TTTFF) - E(TF, TTF)$$

Selecting Attribute

- Measure the entropy of the original set

~~$\frac{3}{5}/\frac{2}{5}$~~

$$E(TTTFF) = -P_T \log_2 P_T - P_F \log_2 P_F = -3/5 \cdot \log_2 3/5 - 2/5 \cdot \log_2 2/5 = 0.97$$

- Measure the entropy of each split

$$E(TT) = -P_T \log_2 P_T - P_F \log_2 P_F = -2/2 \cdot \log_2 2/2 - 0/2 \cdot \log_2 0/2 = 0.0$$

$$E(TFF) = -P_T \log_2 P_T - P_F \log_2 P_F = -1/3 \cdot \log_2 1/3 - 2/3 \cdot \log_2 2/3 = 0.92$$

$E(TT, TFF)$

$$\frac{2}{5} \times 0 + \frac{3}{5} \times 0.92 = 0.55$$

$E(T, TTFF)$

$$\frac{1}{5} \times 0 + \frac{4}{5} \times 1.0 = 0.80$$

$E(TF, TTF)$

$$\frac{2}{5} \times 1.0 + \frac{3}{5} \times 0.92 = 0.952$$

$$E(TF) = -P_T \log_2 P_T - P_F \log_2 P_F = -1/2 \cdot \log_2 1/2 - 1/2 \cdot \log_2 1/2 = 1.0$$

$$E(TTF) = -P_T \log_2 P_T - P_F \log_2 P_F = -2/3 \cdot \log_2 2/3 - 1/3 \cdot \log_2 1/3 = 0.92$$

Selecting Attribute

- Gain: Measure how much an attribute reduce the entropy

$$Gain(S, a) = Entropy(S) - Entropy \text{ of splits by } a$$

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

A=T

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T

A=F

Day	A	B	C	Play
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

B=T

Day	A	B	C	Play
3	F	T	F	T

B=F

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
4	F	F	T	F
5	F	F	F	F

Day	A	B	C	Play
1	T	F	T	T
2	T	F	F	T
3	F	T	F	T
4	F	F	T	F
5	F	F	F	F

C=T

Day	A	B	C	Play
1	T	F	T	T
4	F	F	T	F

C=F

Day	A	B	C	Play
2	T	F	F	T
3	F	T	F	T
5	F	F	F	F

Decrease of Impurity

$$Gain(Play, A) = 0.97 - 0.55 = 0.42$$

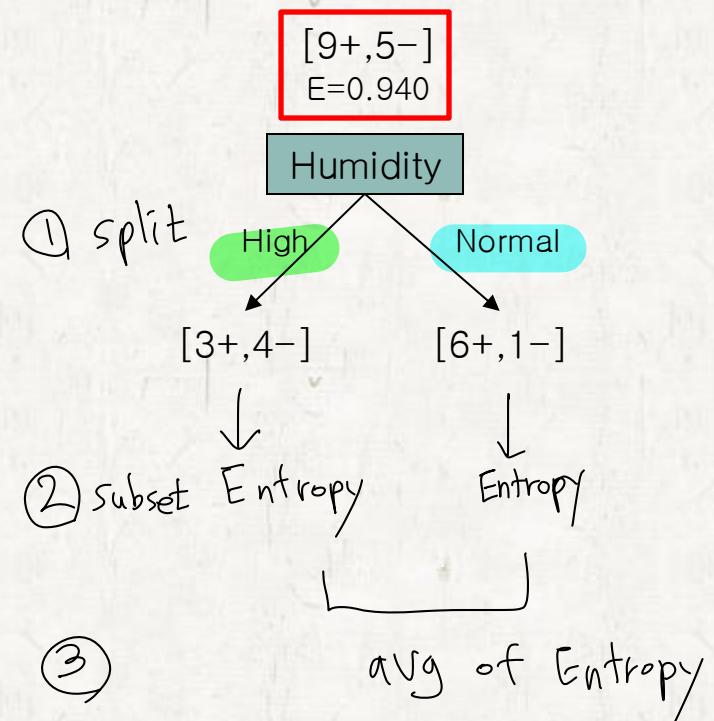
$$Gain(Play, B) = 0.97 - 0.80 = 0.17$$

$$Gain(Play, C) = 0.97 - 0.952 = 0.018$$

Example

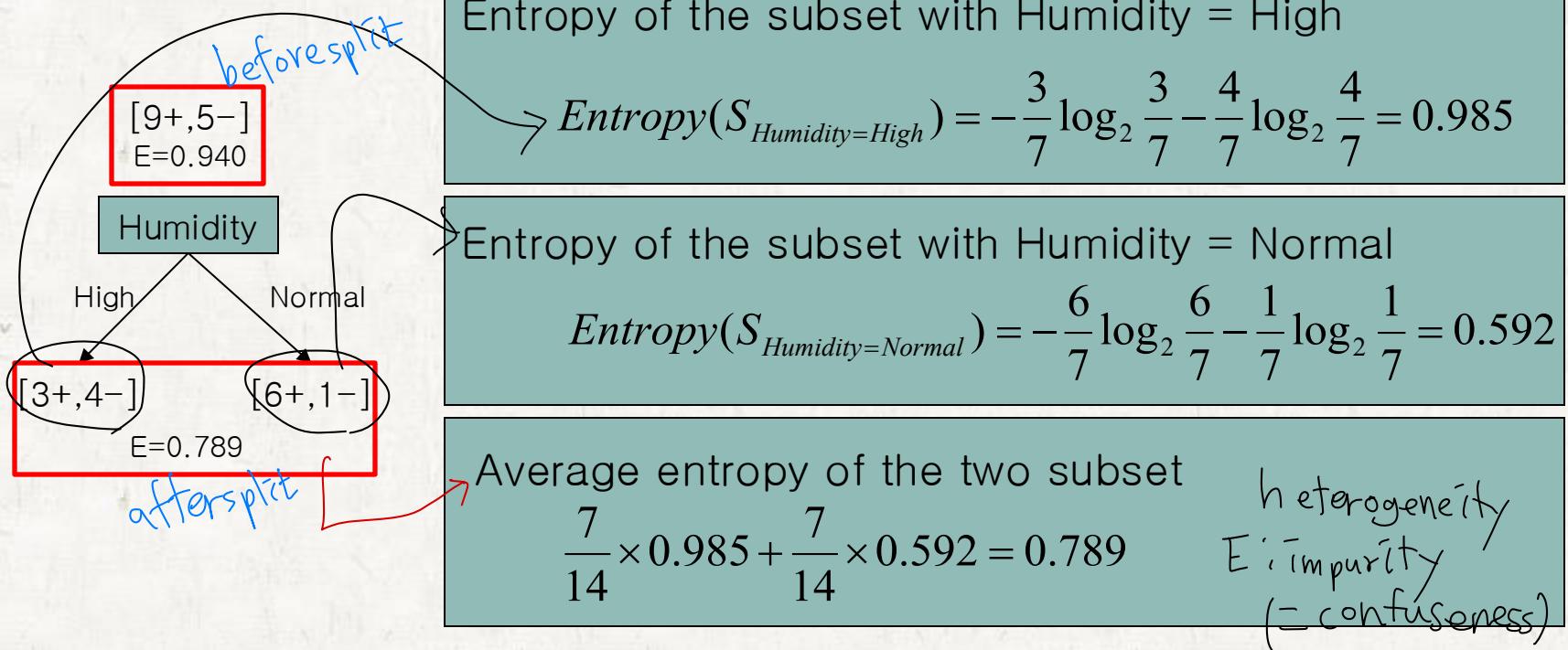
- Tennis playing record again

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No



Example

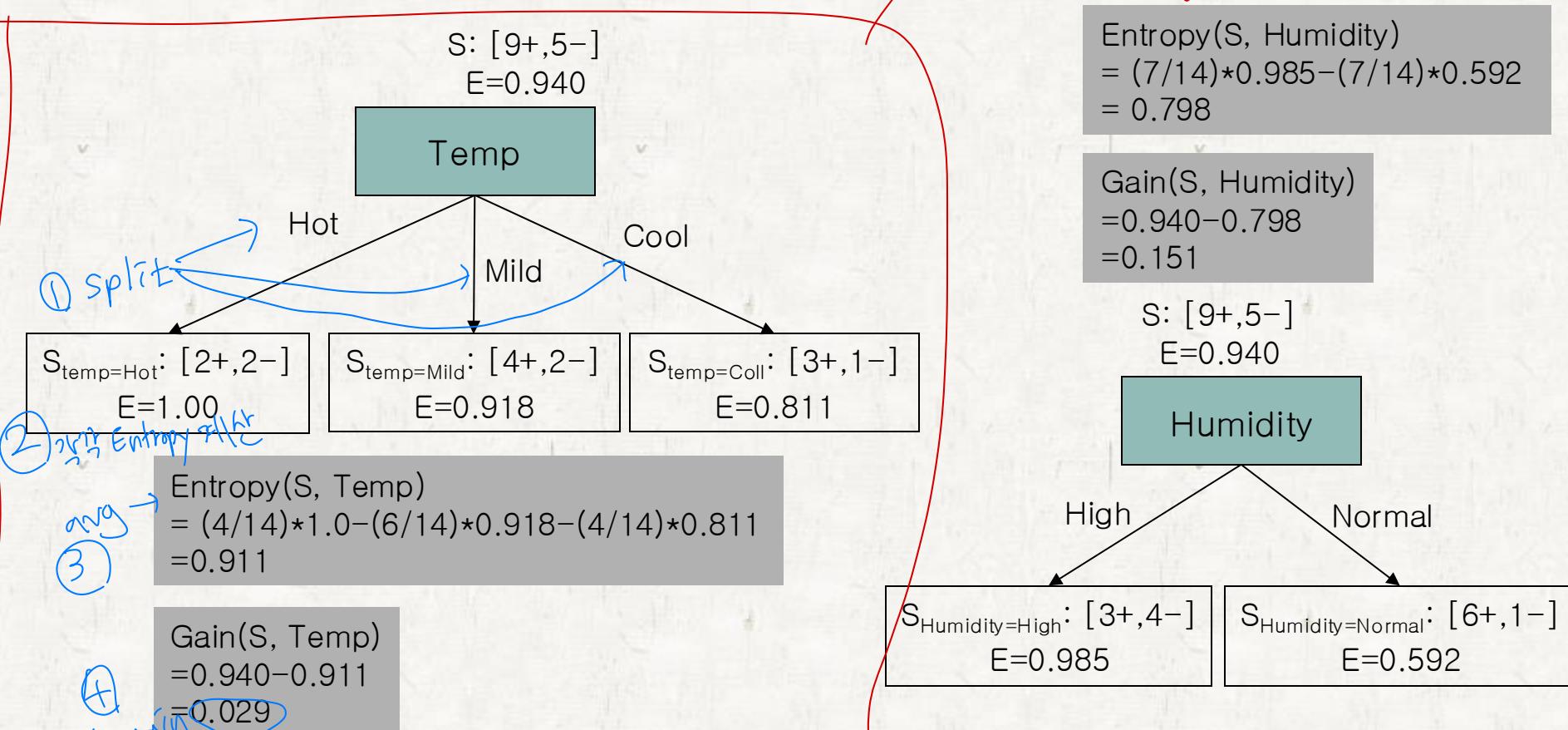
- Reduction of disorder after dividing
 - S is divided into two subsets
 - What is the average entropy of the two subsets ?



Example

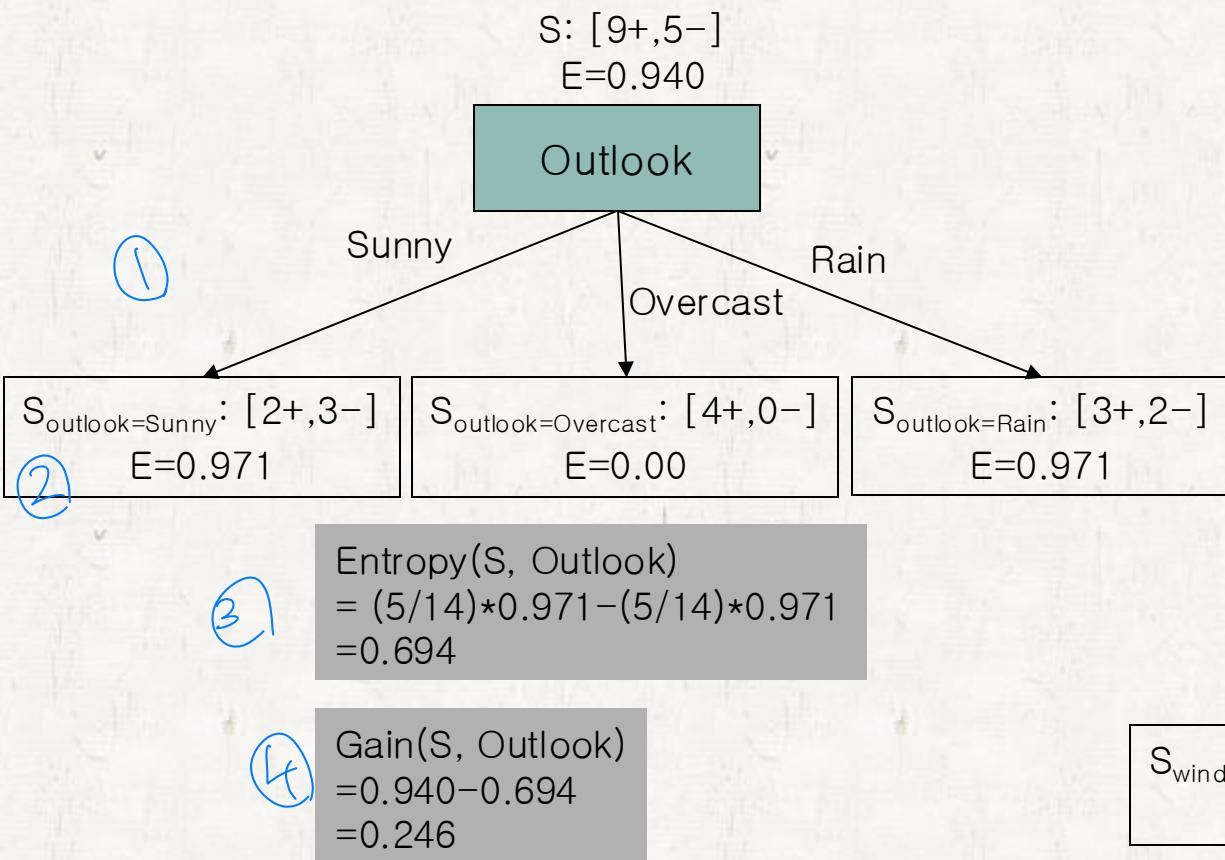
- Which attribute is the best classifier?

ID3 방법!
Base: ID3
 그 외에 C4.5, CART 등..
 방법 존재



Example

- Which attribute is the best classifier?



$$\begin{aligned} \text{Entropy}(S, \text{Wind}) &= (8/14)*0.811 - (6/14)*1.00 \\ &= 0.892 \end{aligned}$$

$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= 0.940 - 0.892 \\ &= 0.048 \end{aligned}$$

S: [9+,5-]
E=0.940

Wind



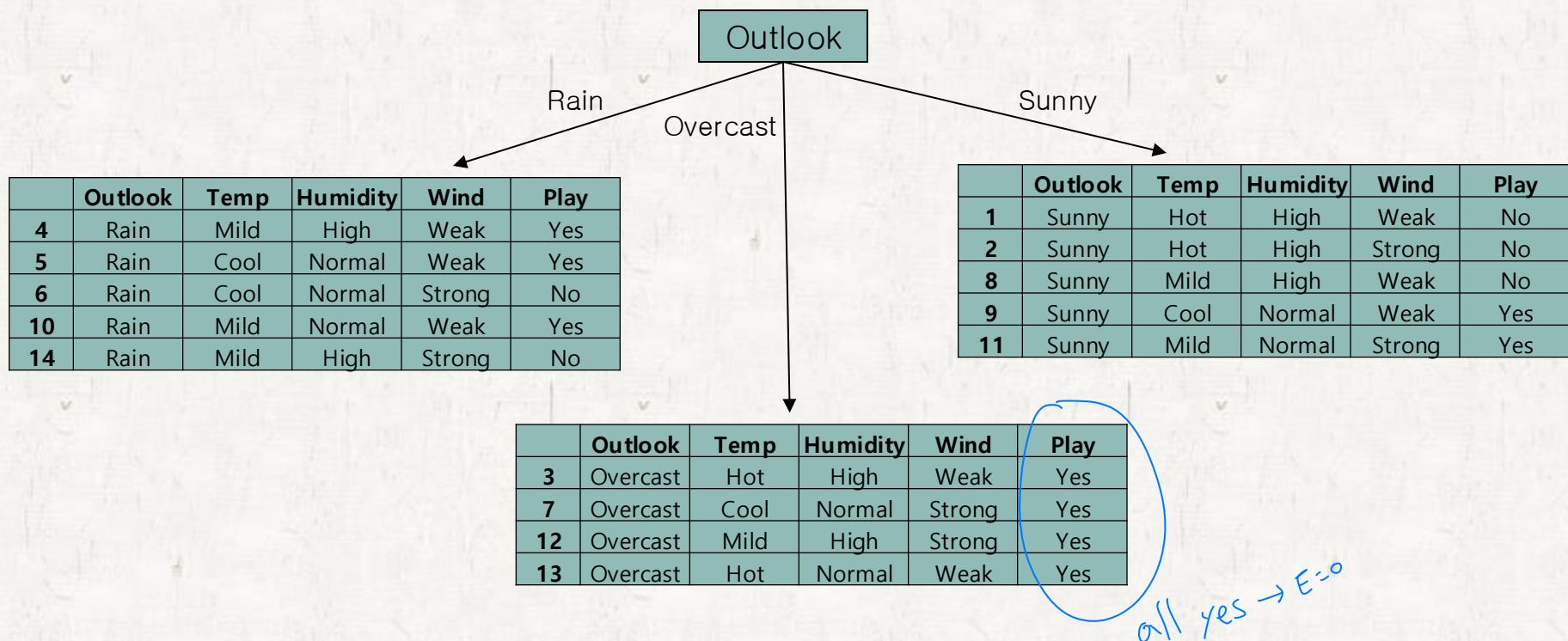
Example

- Which attribute is the best classifier?
 - At root node, we will classify the data with “Outlook”

Gain(S, Outlook) = 0.246
Gain(S, Humidity) = 0.151
Gain(S, Wind) = 0.048
Gain(S, Temp) = 0.029

Example

- Which attribute is the best classifier?
 - At root node, we will classify the data with “Outlook”



Example

- Which attribute is the best classifier?

- At root node, we will classify the data with “Outlook”

Outlook

- Rain
- Overcast
- Sunny

Yes

outlook⁰³
우선 향하고
우선 향하고

E↓⁰³ (best split)
자른 향이며

Yet, not completely separated. What will you do?

Recursively, apply the same procedure !!

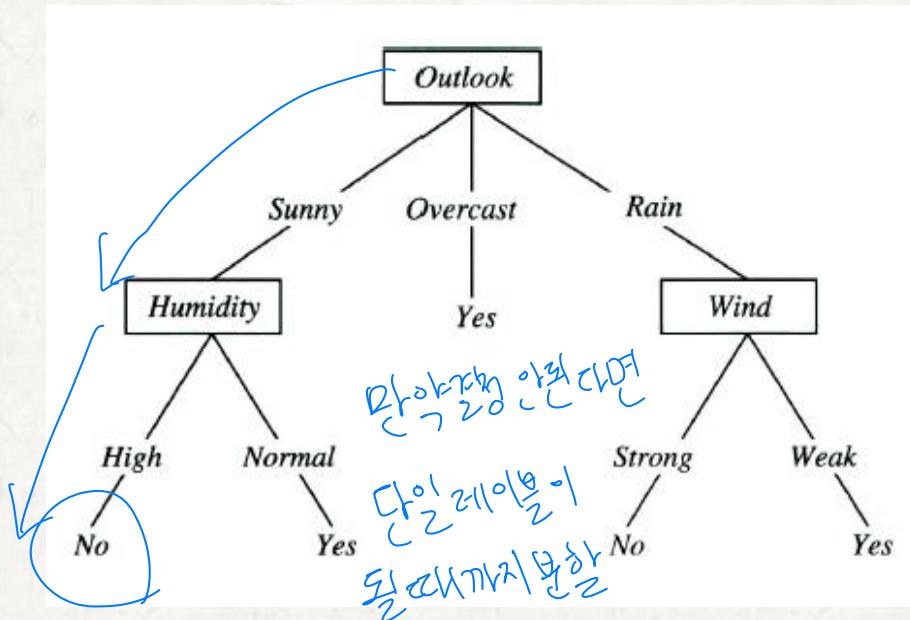
	Outlook	Temp	Humidity	Wind	Play
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
10	Rain	Mild	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes

Example

- Created Decision Tree

→ 노는 예상치 않아
방법 같은 타운 도출하지 않는다



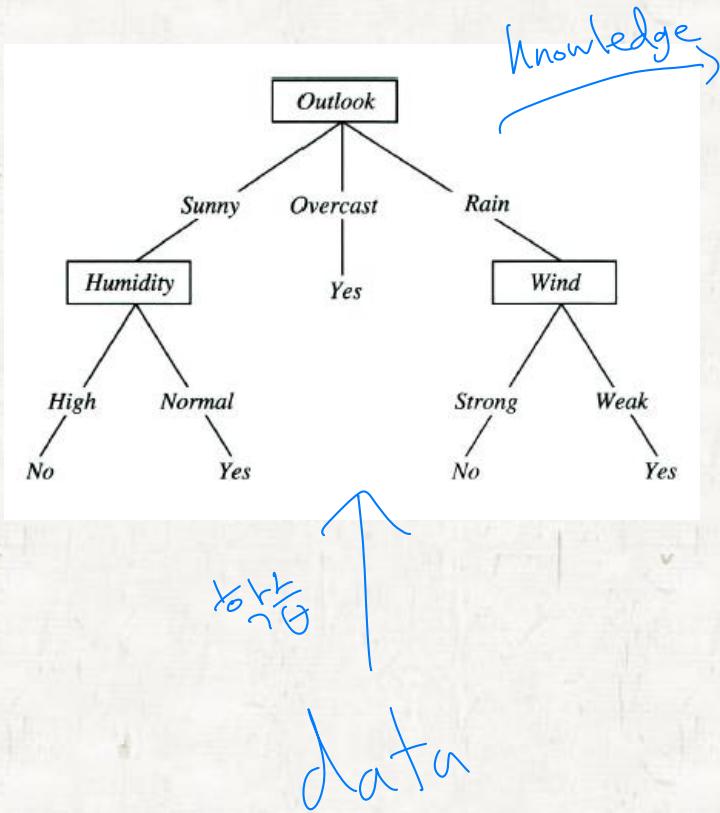
- Prediction

- Sunny, Mild, High, Strong → No

Example

interpretable 理解可能

Rule Generation from Decision Tree



If Outlook = Sunny & Humidity = High
then Tennis = No

If Outlook = Sunny & Humidity = Normal
then Tennis = Yes

If Outlook = Overcast
then Tennis = Yes

If Outlook = Rain & Wind = Strong
then Tennis = No

If Outlook = Sunny & Wind = Weak
then Tennis = Yes

37H Questions

2.
 - There is the unique Decision Tree for a given dataset "No"
 - Which tree is the best?
2.
 - ID3 generates the best Decision Tree "No" ID3는 Best가 아님
3. ID3 guarantees "the best" tree "No"

Questions

- ① choose an Input
- ② split the table (puretable)
- ③ Repeat this

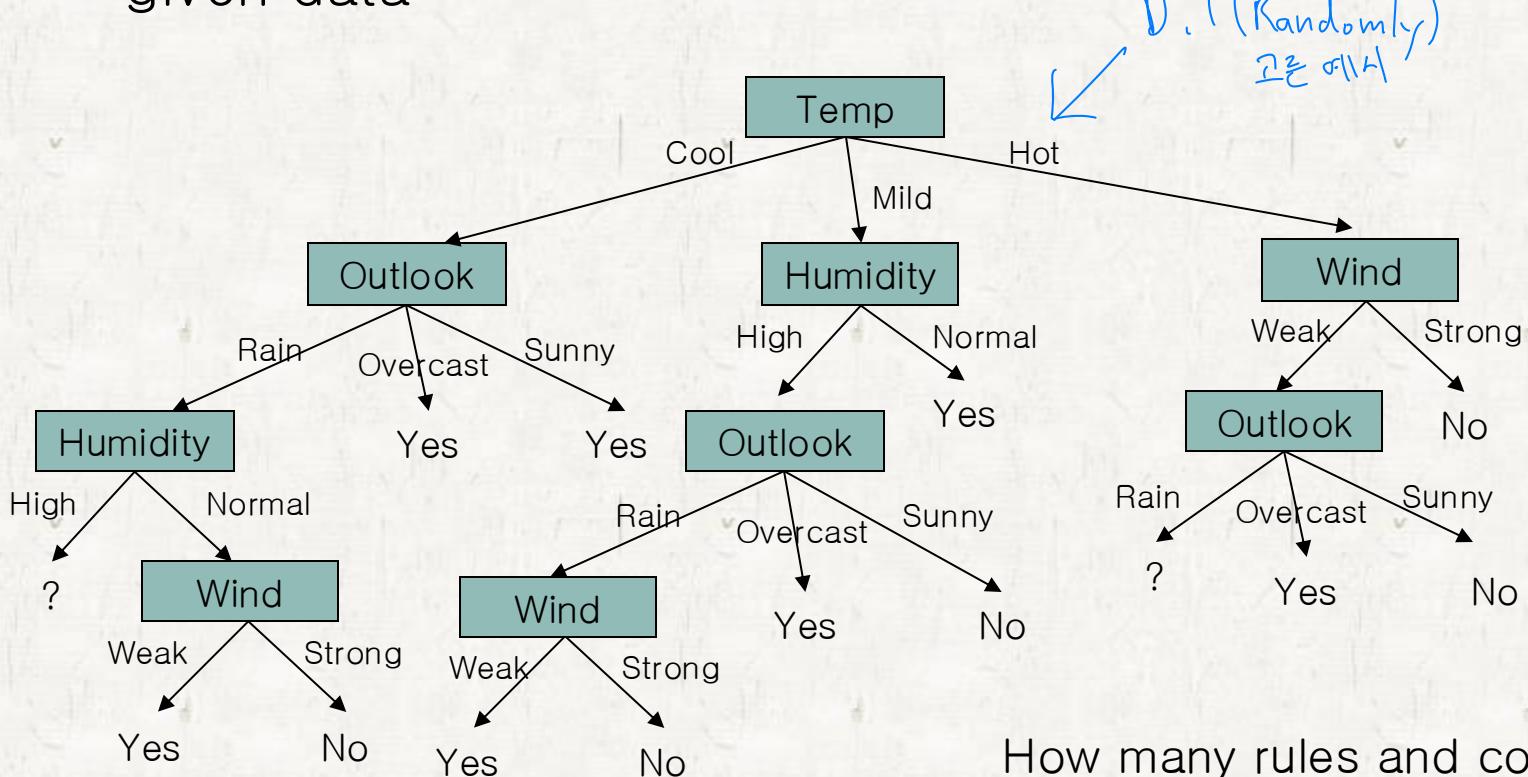
Entropy,

GINI

Randomly

(best among 가능)

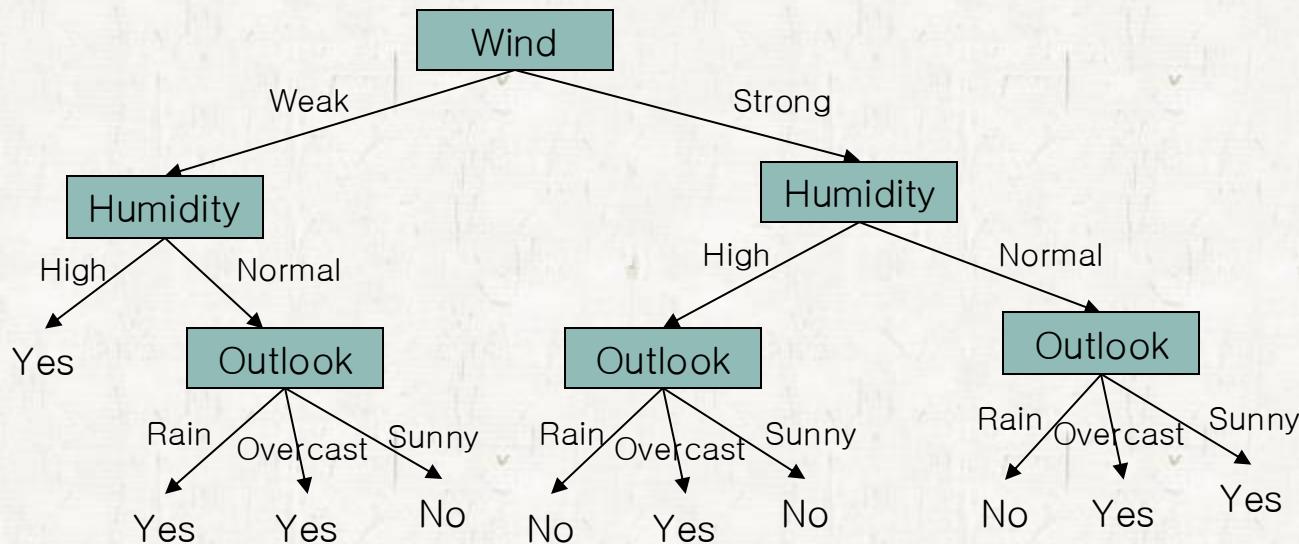
- Decision Tree by ID3 is not the unique decision tree for given data



How many rules and conditions?

Questions

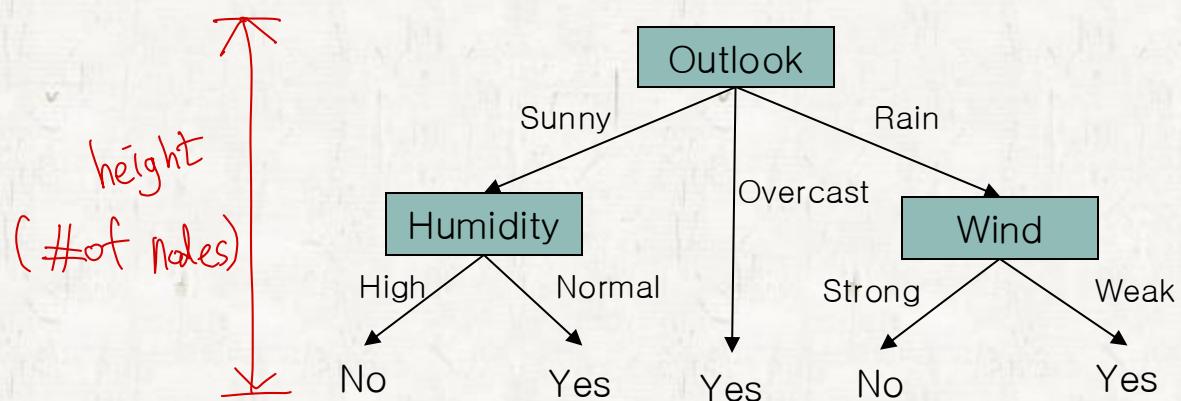
- Decision Tree by ID3 is not unique decision tree for given data



How many rules and conditions?

Questions

- Decision Tree by ID3 is not unique decision tree for given data



D.T (Entropy)

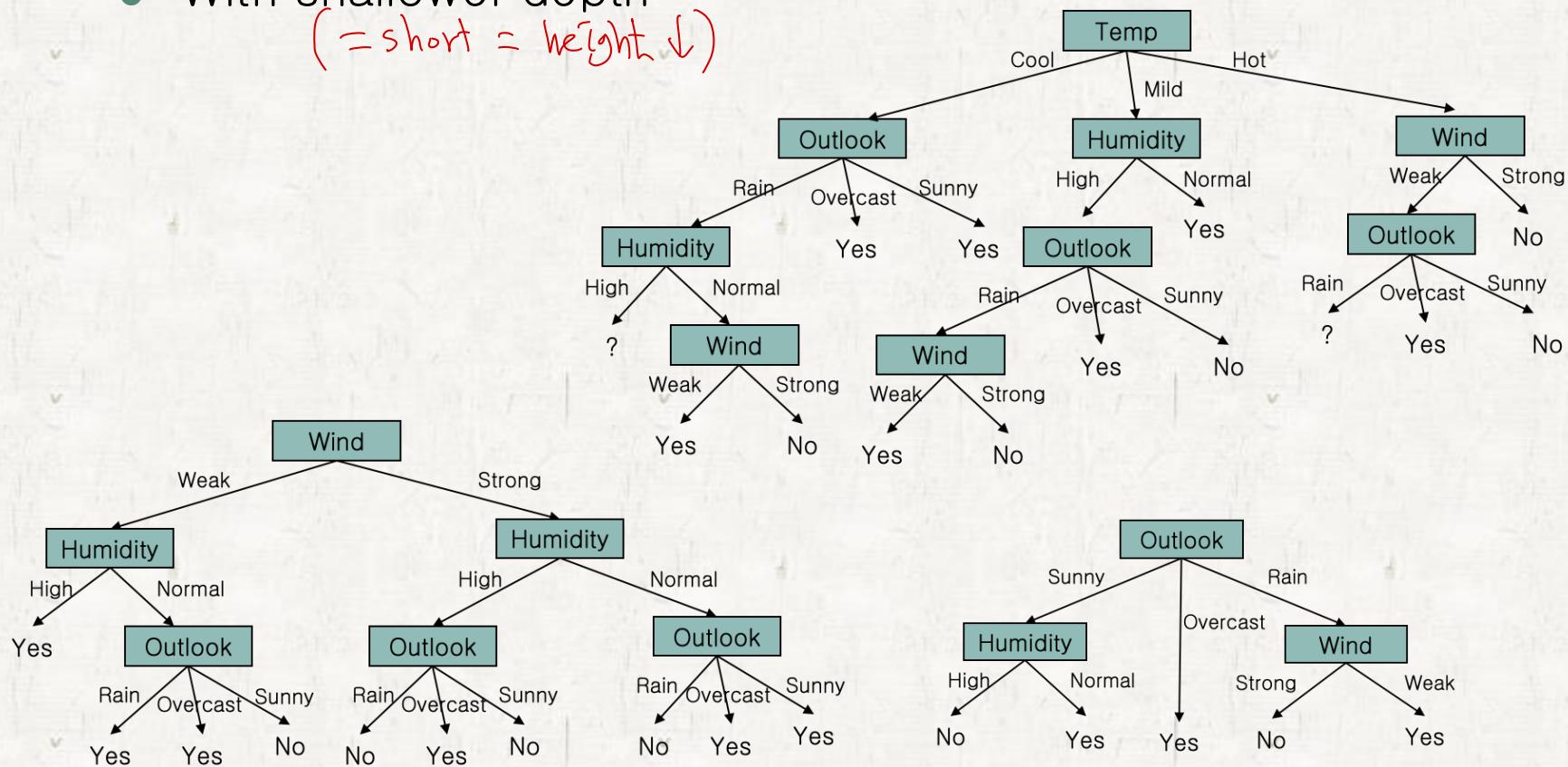
"The simpler is the better"

위험, 복잡한 것은 overfitting 드는
쉽기 때문에 (주어진 데이터에)

How many rules and conditions?

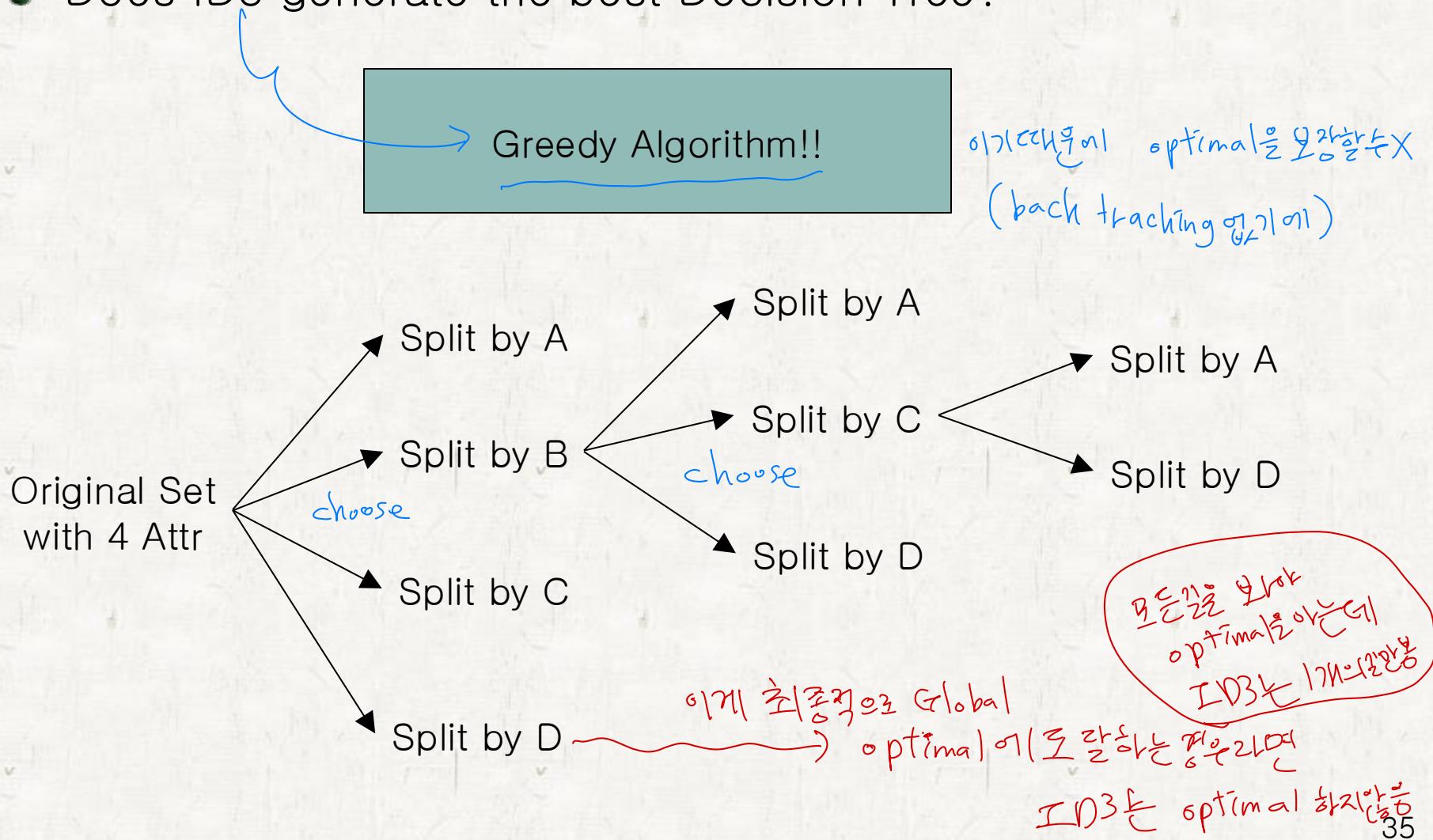
Questions

- Which tree is the best?
*(= simple
= short = height ↓)*
- With fewer nodes
- With shallower depth



Questions

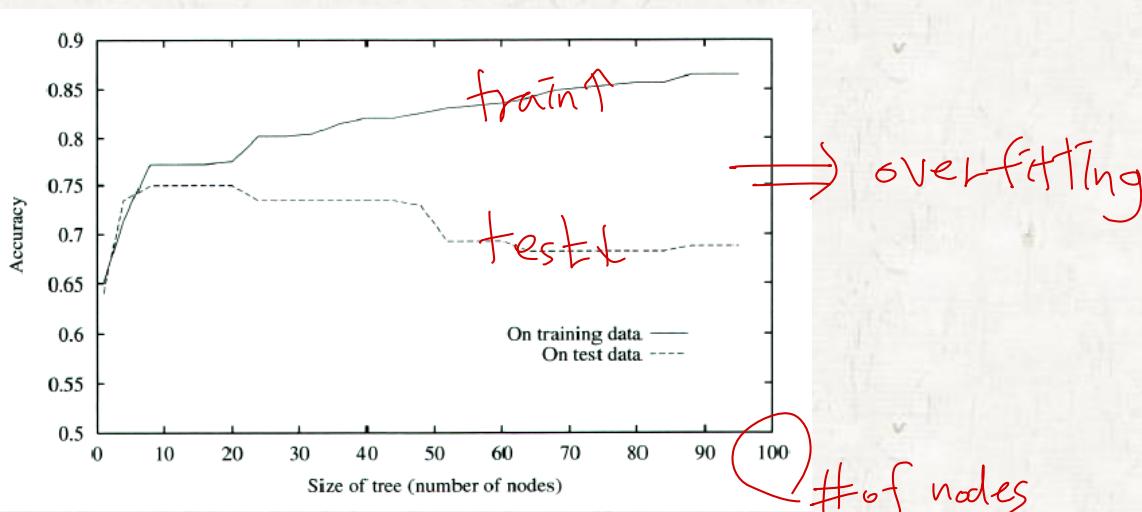
- Does ID3 generate the best Decision Tree?



Discussion

Overfitting

- As # of nodes is increasing, the Decision Tree fits to training data more and more, but possibly overfitted



How to avoid overfitting

- When generating: limit the maximum depth, or do not split nodes with small samples
- After generating: Pruning

Table →

Table의
모든 노드가 학습 데이터
(장점) but overfitting인
(단점)

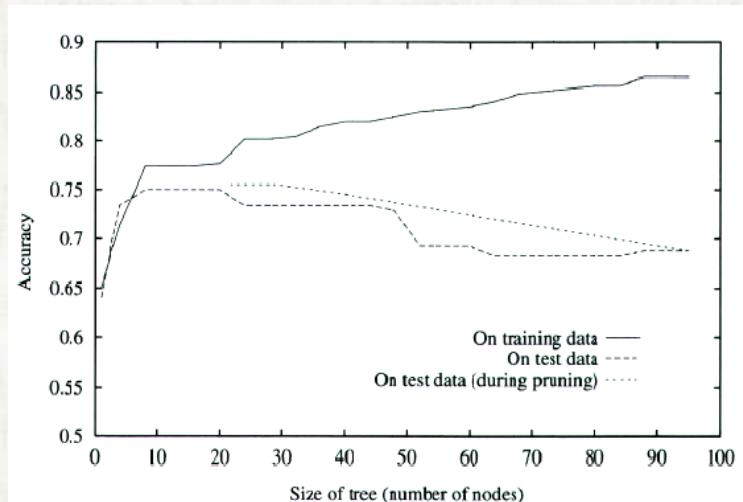
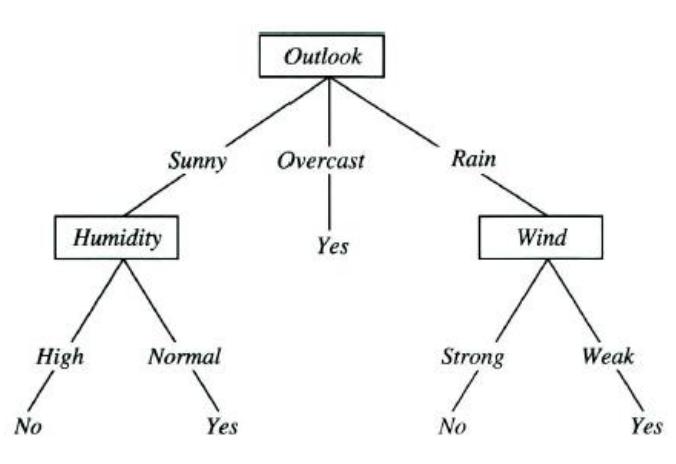


Discussion

Pruning

- Reduced Error Pruning

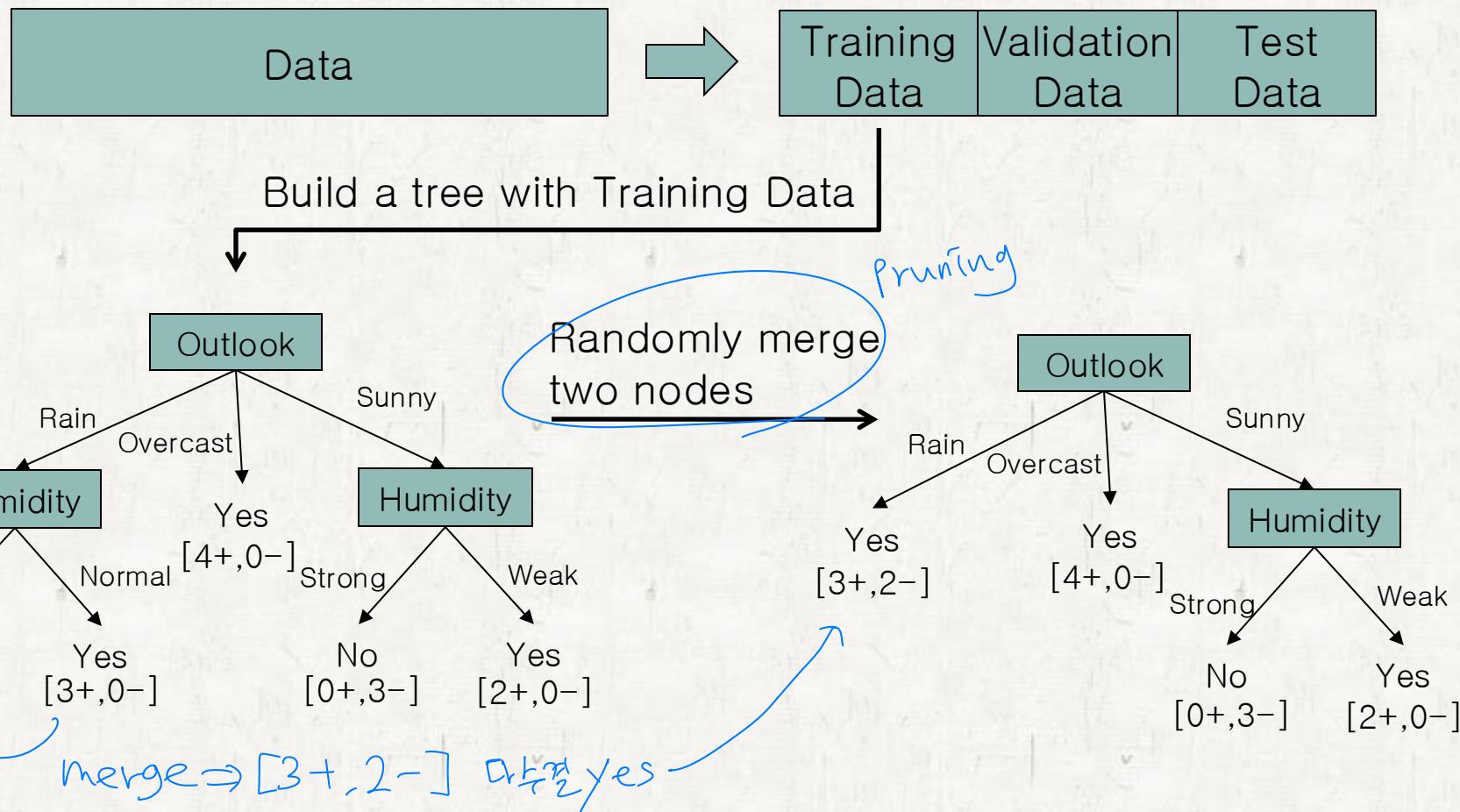
- Divide data into three set : training set, validation set, testing set
- Create a decision tree with training set (allow overfit)
- Prune a leaf node if removing it does not affect the classification of validation set
- Do this until there is no more such nodes



Discussion

Pruning

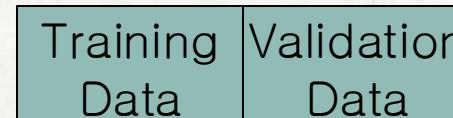
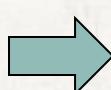
- Reduced Error Pruning



Discussion

Pruning

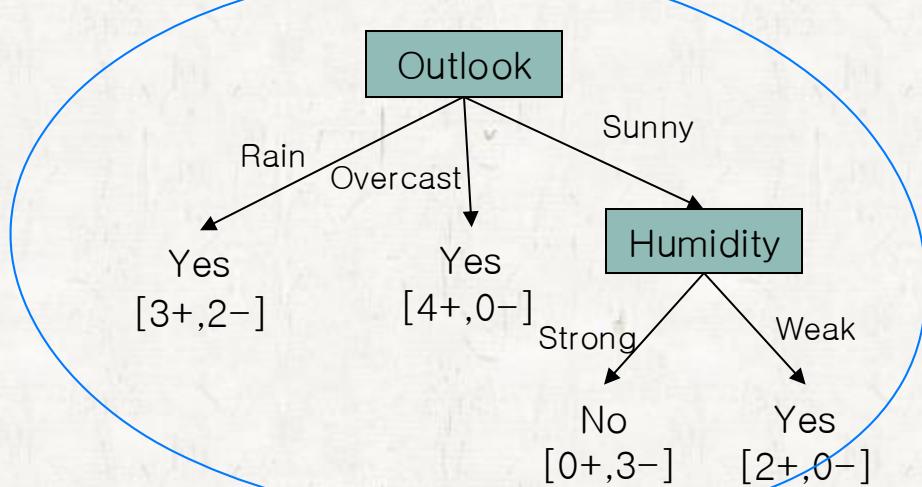
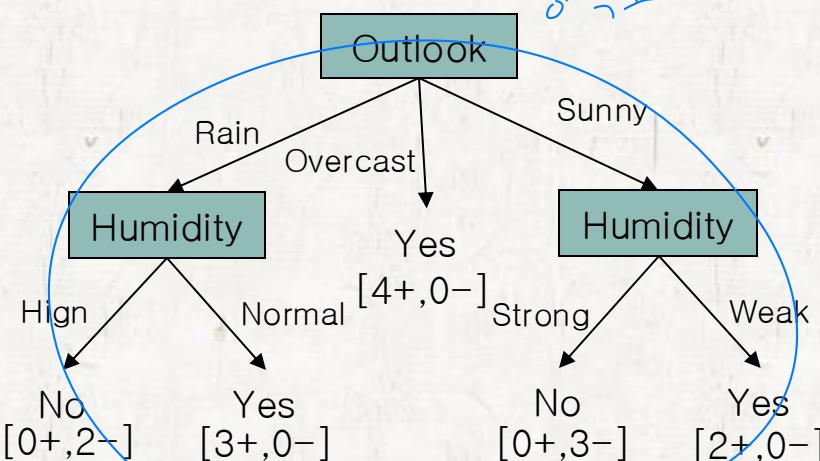
- Reduced Error Pruning



o Tree를 사용할지 결정하는 단계
(성능측정)

Evaluate ~~trees~~ with validation data

정확도측정



pruning 후 정확도가 더 높으면 치적! (이 과정 반복)

Discussion

optimal tree는
불가능하다

Pruning

Post Rule Pruning

또 다른 방법

- Split data into training set and validation set
- Build a consistent decision tree (allow overfit)
= full, complete
- Convert the tree into a set of "if-then" rules
- For each rule, remove any attributes if estimated accuracy of validation set is not reduced

↓
조건

If (Outlook = Sunny) & (Humidity = High) then Tennis = No
If (Outlook = Sunny) & (Humidity = Normal) then Tennis = Yes
If (Outlook = Overcast) then Tennis = Yes
If (Outlook = Rain) & (Wind = Strong) then Tennis = No
If (Outlook = Sunny) & (Wind = Weak) then Tennis = Yes

⇒ 9개 조건 중에 1개를 지우면 그 후 정확도를 측정함

지우기전 / 후 정확도 중 높은걸 고름