

ID3

C4.5, CART

셋의 공통 알고리즘

이 고정의 차이가  
셋을 구분시킴

- ① choose an input
  - ② split table
  - ③ repeat these steps
- Until pure subtable

Categorical ✓

Real ✓

Entropy ✓

Gini ✓

Gain ✓

GainInfo<sup>v</sup>

모두 동일이라 CNT(F) 사용 가능하다

Categorical

C4.5

ID3

⇒ 범주형 데이터

숫자형데이터

## Dataset with Continuous-Valued Attributes

OUTLOOK	TEMP.	HUMIDITY	WINDY	PLAY
sunny	85	85	false	No
sunny	80	90	true	No
overcast	83	78	false	Yes
rain	70	96	false	Yes
rain	68	80	false	Yes
rain	65	70	true	No
overcast	64	65	true	Yes
sunny	72	95	false	No
sunny	69	70	false	Yes
rain	75	80	false	Yes
sunny	75	70	true	Yes
overcast	72	90	true	Yes
overcast	81	75	false	Yes
rain	71	80	true	No

(X, Y) X: 입력 Y: 출력

① X가 숫자형

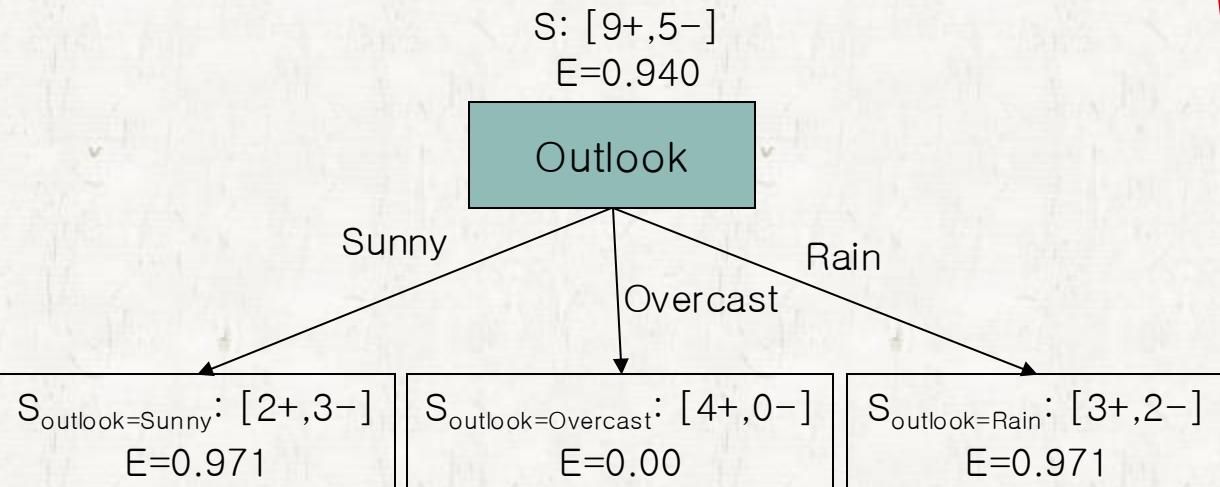
② Y가 숫자형

# C4.5

Categorical  
자료형

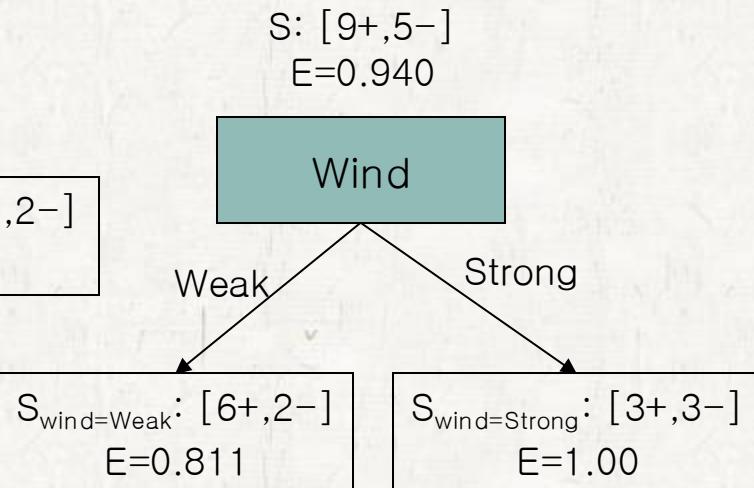
1. Evaluate the entropy of the given table
  2. Split the table with an input
  3. Evaluate the entropy of each sub-table
  4. Evaluate the average of the entropies
  5. Evaluate the Gain: (1)–(4)
- \* Repeat this for all inputs and choose the best

## Categorical Attribute: Split as we did



$$\begin{aligned} \text{Gain}(S, \text{Outlook}) \\ = 0.940 - (5/14)*0.971 - (5/14)*0.971 \\ = 0.246 \end{aligned}$$

- ① Choose an input
- Try all inputs
  - ① Split with an input
  - ② Entropy
  - ③ Gain
  - choose the best
- ② Split table



$$\begin{aligned} \text{Gain}(S, \text{Wind}) \\ = 0.940 - (8/14)*0.811 - (6/14)*1.00 \\ = 0.048 \end{aligned}$$

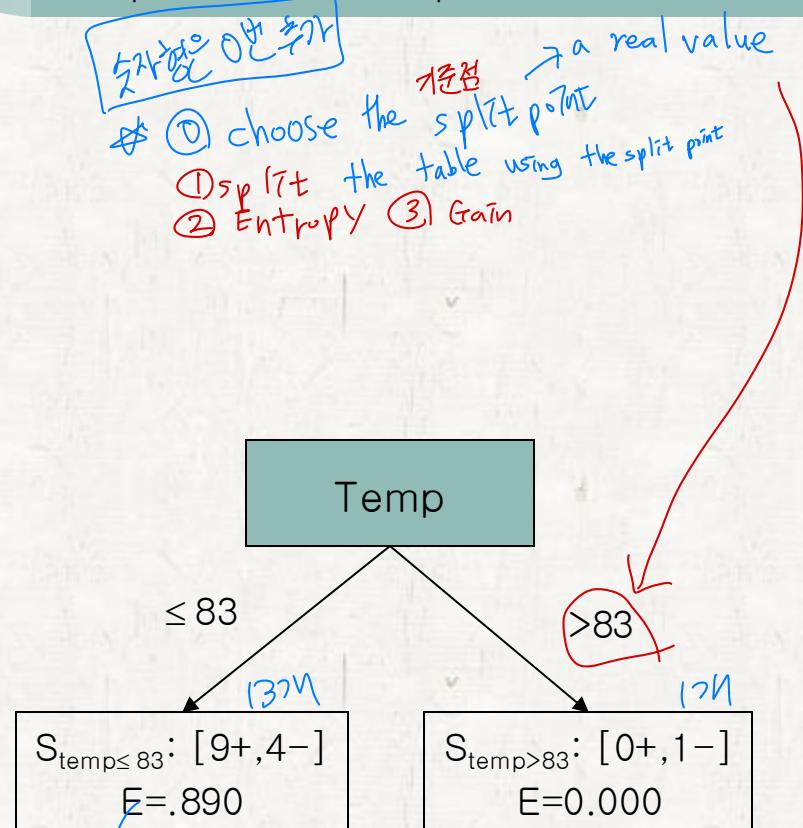
# C4.5

1. Evaluate the entropy of the given table
  2. Split the table with an input
  3. Evaluate the entropy of each sub-table
  4. Evaluate the average of the entropies
  5. Evaluate the Gain: (1)–(4)
- \* Repeat this for all inputs and choose the best

## Continuous-Valued Attributes

- Binary Split

OUTLOOK	TEMP.	HUMIDITY	WINDY	PLAY
sunny	85	85	false	No
sunny	80	90	true	No
overcast	83	78	false	Yes
rain	70	96	false	Yes
rain	68	80	false	Yes
rain	65	70	true	No
overcast	64	65	true	Yes
sunny	72	95	false	No
sunny	69	70	false	Yes
rain	75	80	false	Yes
sunny	75	70	true	Yes
overcast	72	90	true	Yes
overcast	81	75	false	Yes
rain	71	80	true	No



$$\text{Gain} = 0.940 - \frac{13}{14} \times 0.890$$

How to choose the split point?  
 $\text{Temp} \rightarrow 81, 3, 5, 15$

$$\text{SplitInfo} = -\frac{13}{14} \log \frac{13}{14} - \frac{1}{14} \log \frac{1}{14}$$

Entropy → E는 두 집합의 Entropy의 평균으로 정함

# C4.5

Continuous input  
k(210)log

1. Evaluate the entropy of the given table
  - 2-1. Choose the split point
  - 2-2. Split the table with the split point
  3. Evaluate the entropy of each sub-table
  4. Evaluate the average of the entropies
  5. Evaluate the Gain: (1)-(4)
- \* Repeat this for all inputs and choose the best

## Continuous-Valued Attributes

- How to choose the split point?

OUTLOOK	TEMP.	HUMIDITY	WINDY	PLAY
sunny	85	85	false	No
sunny	80	90	true	No
overcast	83	78	false	Yes
rain	70	96	false	Yes
rain	68	80	false	Yes
rain	65	70	true	No
overcast	64	65	true	Yes
sunny	72	95	false	No
sunny	69	70	false	Yes
rain	75	80	false	Yes
sunny	75	70	true	Yes
overcast	72	90	true	Yes
overcast	81	75	false	Yes
rain	71	80	true	No

All values in the dataset is  
the candidates of split point

- ① Evaluate all split point  
② choose the best one

# C4.5

1. Evaluate the entropy of the given table
  - 2-1. Choose the split point
  - 2-2. Split the table with the split point
  3. Evaluate the entropy of each sub-table
  4. Evaluate the average of the entropies
  5. Evaluate the Gain: (1)–(4)
- \* Repeat this for all inputs and choose the best

## How to choose the split point

- Step 1: Try all the candidate split points

OUTLOOK	TEMP.	HUMIDITY	WINDY	PLAY
overcast	64	65	true	Yes
rain	65	70	true	No
rain	68	80	false	Yes
sunny	69	70	false	Yes
rain	70	96	false	Yes
rain	71	80	true	No
sunny	72	95	false	No
overcast	72	90	true	Yes
rain	75	80	false	Yes
sunny	75	70	true	Yes
sunny	80	90	true	No
overcast	81	75	false	Yes
overcast	83	78	false	Yes
sunny	85	85	false	No

- 
- (1)  $\leq 64$  : [1+, 0-],  $> 64$ : [8+, 5-]
  - (2)  $\leq 65$  : [1+, 1-],  $> 65$ : [8+, 4-]
  - (3)  $\leq 68$  : [2+, 1-],  $> 68$ : [7+, 4-]
  - (4)  $\leq 69$  : [3+, 1-],  $> 69$ : [6+, 4-]
  - (5)  $\leq 70$  : [4+, 1-],  $> 70$ : [5+, 4-]
  - (6)  $\leq 71$  : [4+, 2-],  $> 71$ : [5+, 3-]
  - (7)  $\leq 72$  : [5+, 3-],  $> 72$ : [4+, 2-]
  - (8)  $\leq 75$  : [7+, 3-],  $> 75$ : [2+, 2-]
  - (9)  $\leq 80$  : [7+, 4-],  $> 80$ : [2+, 1-]
  - (10)  $\leq 81$  : [8+, 4-],  $> 81$ : [1+, 1-]
  - (11)  $\leq 83$  : [9+, 4-],  $> 83$ : [0+, 1-]
  - (12)  $\leq 85$  : [9+, 5-],  $> 85$ : [0+, 0-]

# C4.5

1. Evaluate the entropy of the given table
  - 2-1. Choose the split point
  - 2-2. Split the table with the split point
  3. Evaluate the entropy of each sub-table
  4. Evaluate the average of the entropies
  5. Evaluate the Gain: (1)-(4)
- \* Repeat this for all inputs and choose the best

## How to choose the split point

- Step2: Evaluate split qualities of them

- (1)  $\leq 64$  : [1+, 0-],  $>64$ : [8+, 5-]  $\rightarrow$  Entropy = 0.893
- (2)  $\leq 65$  : [1+, 1-],  $>65$ : [8+, 4-]  $\rightarrow$  Entropy = 0.930
- (3)  $\leq 68$  : [2+, 1-],  $>68$ : [7+, 4-]  $\rightarrow$  Entropy = 0.940
- (4)  $\leq 69$  : [3+, 1-],  $>69$ : [6+, 4-]  $\rightarrow$  Entropy = 0.925
- (5)  $\leq 70$  : [4+, 1-],  $>70$ : [5+, 4-]  $\rightarrow$  Entropy = 0.895
- (6)  $\leq 71$  : [4+, 2-],  $>71$ : [5+, 3-]  $\rightarrow$  Entropy = 0.939
- (7)  $\leq 72$  : [5+, 3-],  $>72$ : [4+, 2-]  $\rightarrow$  Entropy = 0.939
- (8)  $\leq 75$  : [7+, 3-],  $>75$ : [2+, 2-]  $\rightarrow$  Entropy = 0.915
- (9)  $\leq 80$  : [7+, 4-],  $>80$ : [2+, 1-]  $\rightarrow$  Entropy = 0.940
- (10)  $\leq 81$  : [8+, 4-],  $>81$ : [1+, 1-]  $\rightarrow$  Entropy = 0.930
- (11)  $\leq 83$  : [9+, 4-],  $>83$ : [0+, 1-]  $\rightarrow$  Entropy = 0.827
- (12)  $\leq 85$  : [9+, 5-],  $>85$ : [0+, 0-]  $\rightarrow$  Entropy = 0.940

avg Entropy

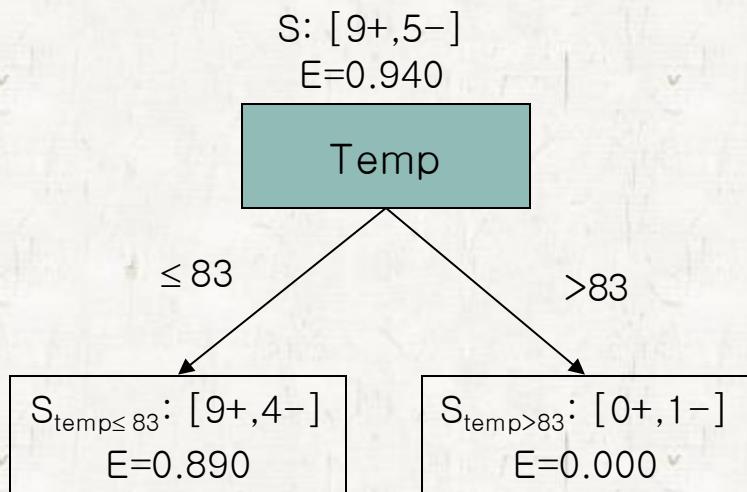
Entropy = confuseness

min

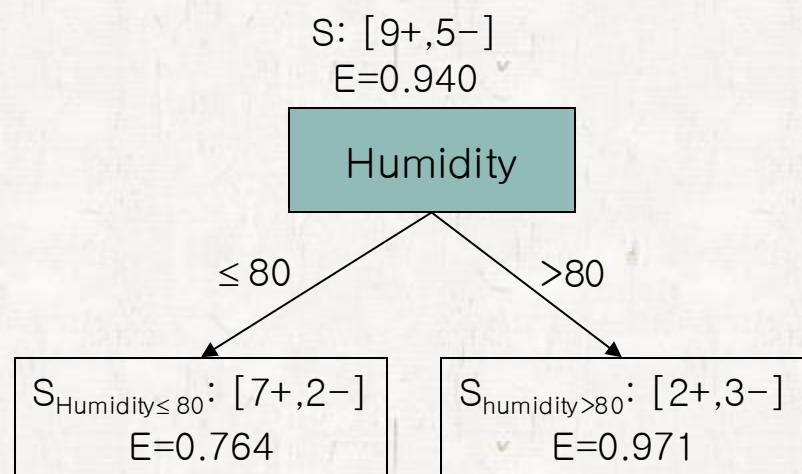
0.827

# C4.5

- Dataset with Continuous-Valued Attributes



$$\begin{aligned} \text{Gain}(S, \text{Temp}) &= 0.940 - 0.827 \\ &= 0.113 \end{aligned}$$



$$\begin{aligned} \text{Gain}(S, \text{Humidity}) &= 0.940 - 0.838 \\ &= 0.102 \end{aligned}$$

$\text{Gain} \leftarrow$   
 Large Gain ↑  
 Large Tree ↓  
 overfitted

ID3와 동일하게

Gain이 가장 높은 것 부터 고름

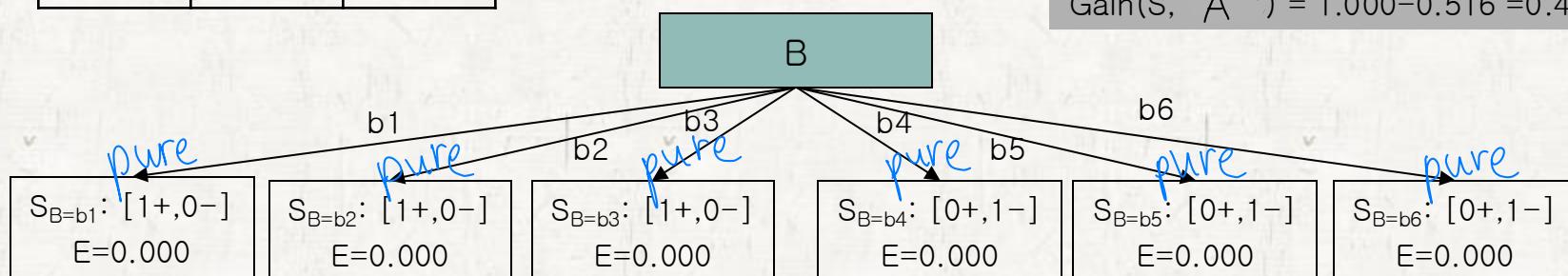
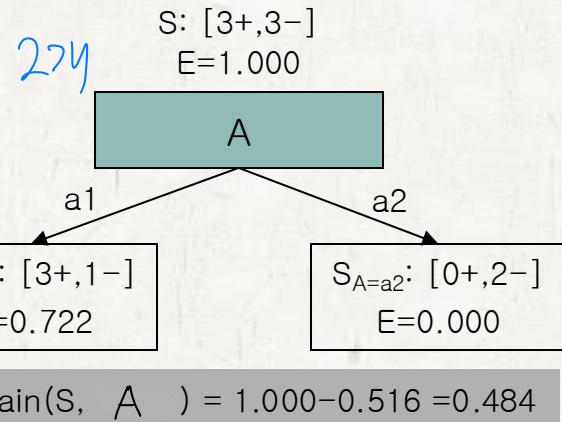
## C4.5

### Alternative Measures for Selecting Attributes (1)

A	B	PLAY
a1	b1	Yes
a1	b2	Yes
a1	b3	Yes
a1	b4	No
a2	b5	No
a2	b6	No

branching factor

f7A  
S: [3+,3-]  
E=1.000



Which one is better? Gain does not consider the branching factor

(Gain을 사용해보면  
주목점! (오버피팅)  
가능성)

# C4.5

- Alternative Measures for Selecting Attributes (2)
  - Considering not only Gain but also branching factor
- SplitInfo: Measuring branching factor
  - Entropy of the attribute

A	B	PLAY
a1	b1	Yes
a1	b2	Yes
a1	b3	Yes
a1	b4	No
a2	b5	No
a2	b6	No

$$A = \{a1, a1, a1, a1, a2, a2\}$$

$$\text{SplitInfo}(\text{set}, A) = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.918$$

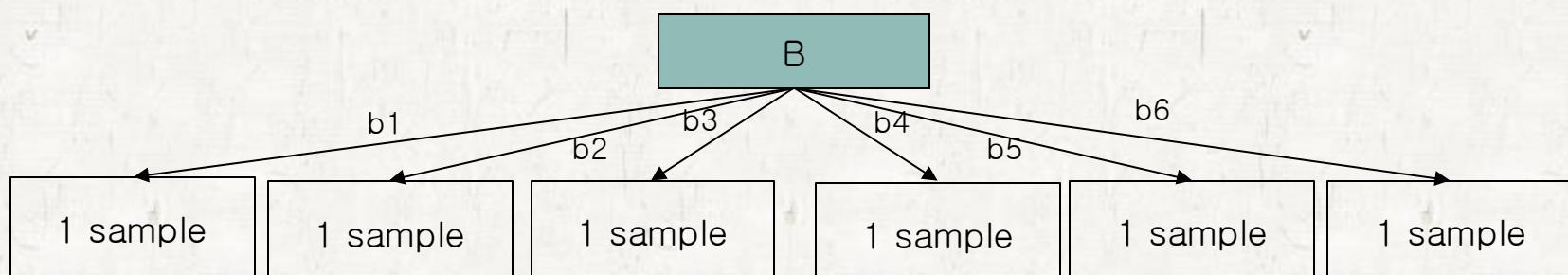
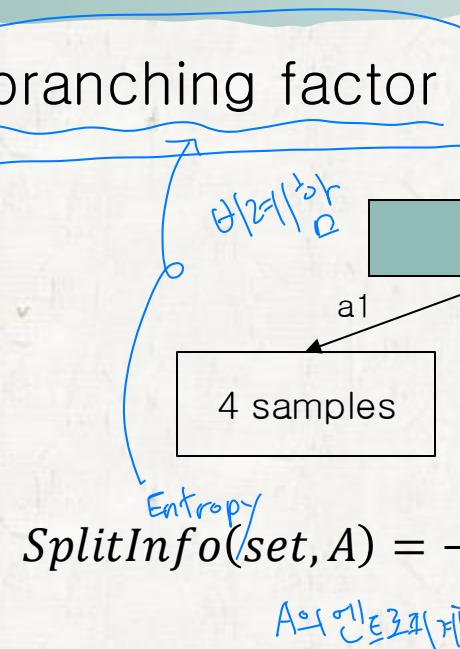
$$B = \{b1, b2, b3, b4, b5, b6\}$$

$$\text{SplitInfo}(\text{set}, B) = 6 \left( -\frac{1}{6} \log_2 \frac{1}{6} \right) = 2.585$$

# C4.5

- SplitInfo: Measuring branching factor

A	B	PLAY
a1	b1	Yes
a1	b2	Yes
a1	b3	Yes
a1	b4	No
a2	b5	No
a2	b6	No



$$SplitInfo(set, B) = 6 \left( -\frac{1}{6} \log_2 \frac{1}{6} \right) = 2.585$$

## C4.5

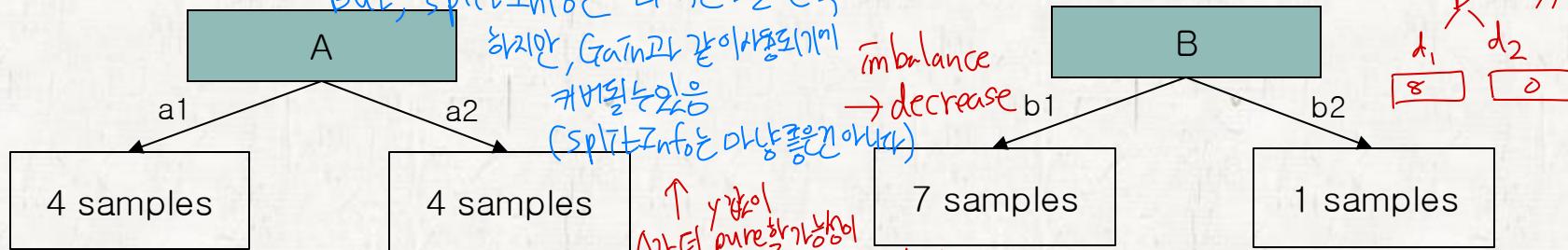
A, C

$\Rightarrow \text{splitInfo} \propto B \cdot F$

$A, B \Rightarrow \text{splitInfo} \propto \text{balance}$

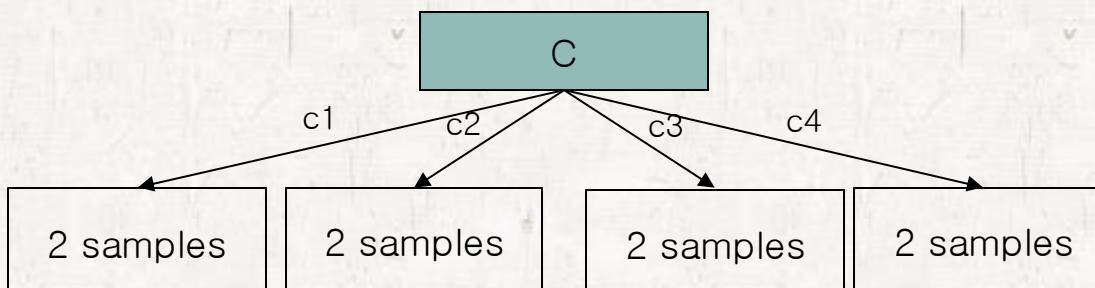
### SplitInfo: Measuring branching factor

D.T가 장애인 트리보다 더 shorter 해지기엔 A를 선택  
But, splitInfo는 더 작은 B를 선택



$$\text{SplitInfo}(\text{set}, A) = -\frac{4}{8} \log_2 \frac{4}{8} - \frac{4}{8} \log_2 \frac{4}{8} = 1.0$$

$$\text{SplitInfo}(\text{set}, B) = -\frac{7}{8} \log_2 \frac{7}{8} - \frac{1}{8} \log_2 \frac{1}{8} = 0.544$$



$$\text{SplitInfo}(\text{set}, C) = -\frac{2}{8} \log_2 \frac{2}{8} - \frac{2}{8} \log_2 \frac{2}{8} - \frac{2}{8} \log_2 \frac{2}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 2$$

## C4.5

usefulness  
of an Input  $\propto \frac{\text{Gain}}{\text{B.F}}$

- Alternative Measures for Selecting Attributes (3)

- Considering not only Gain but also branching factor

$$GainInfo(set, A) = \frac{Gain(set, A)}{SplitInfo(set, A)} \approx \text{measure of branching factor}$$

A	B	PLAY
a1	b1	Yes
a1	b2	Yes
a1	b3	Yes
a1	b4	No
a2	b5	No
a2	b6	No

$$GainInfo(set, A) = \frac{0.484}{0.918} = 0.527$$

$$b.f = 2$$

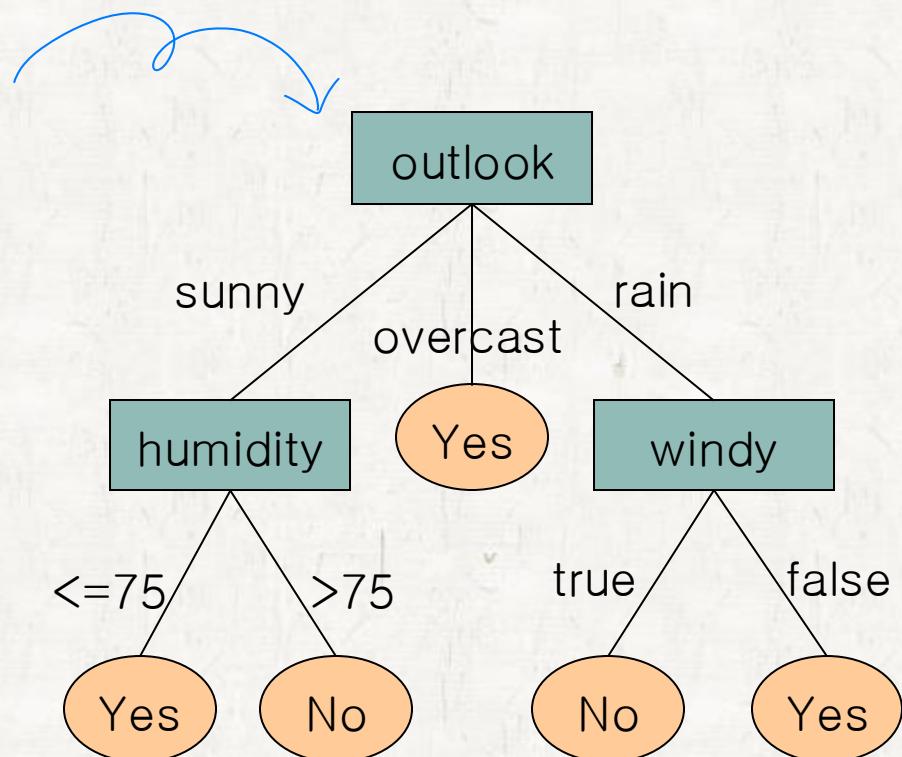
$$GainInfo(set, B) = \frac{1.000}{2.585} = 0.387$$

$$b.f = 6$$

# C4.5

- Incorporating Continuous-Valued Attributes (4)

OUTLOOK	TEMP.	HUMIDITY	WINDY	PLAY
sunny	85	85	false	No
sunny	80	90	true	No
overcast	83	78	false	Yes
rain	70	96	false	Yes
rain	68	80	false	Yes
rain	65	70	true	No
overcast	64	65	true	Yes
sunny	72	95	false	No
sunny	69	70	false	Yes
rain	75	80	false	Yes
sunny	75	70	true	Yes
overcast	72	90	true	Yes
overcast	81	75	false	Yes
rain	71	80	true	No



# CART

- CART developed by Breiman Friedman Olsen and Stone: “Classification and Regression Trees”
- Gini Index instead of Entropy
  - All steps are very similar to C4.5
- Binary Split
  - Eg: Tuesday vs. Not Tuesday
  - Real-Valued Output (Regression Tree)

	ID3	C4.5	CART
Input	Categorical	C, R	$C, R \leftarrow C_{4.5}^{Real\ value}$
Output	C	C	C, R
Split	2, 3, 4..	2, 3, 4... Split Info $\downarrow$	Binary (2)
Measure	Entropy [Gain]	Entropy [GainInfo]	Gini (Index)

이거 3 줄짜고 끝나면

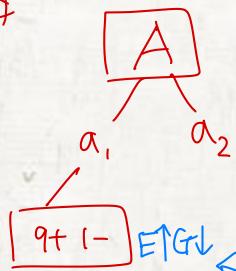
# CART

$$E(\bar{I}) = -\frac{5}{9} \log_2 \frac{5}{9} - \frac{4}{9} \log_2 \frac{4}{9}$$

$$\begin{aligned} Gini(\bar{I}) &= 1 - (p_1^2 + p_0^2) \\ &= 1 - \left( \left(\frac{5}{9}\right)^2 + \left(\frac{4}{9}\right)^2 \right) \end{aligned}$$

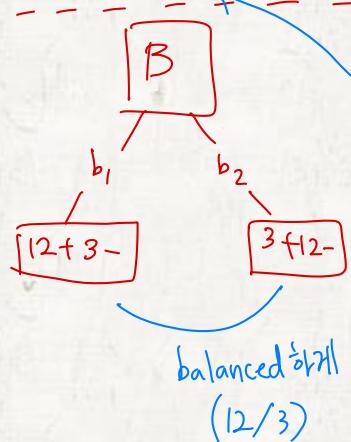
## Gini Index

ex)

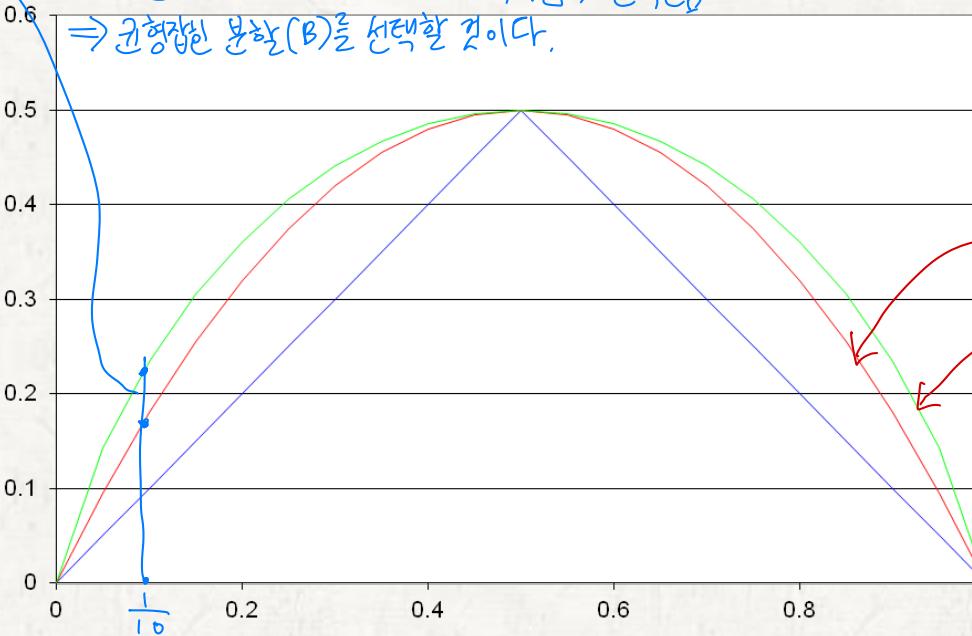


기준레이저고정하고  
한쪽 pure한쪽 다른한쪽 상대적으로 더  
해석

$$i(p) = \sum_{i \neq j} p_i p_j = 1 - \sum_j p_j^2$$

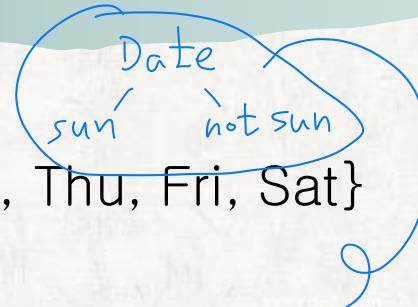


$$i(1111110000) = \frac{5}{9} \cdot \frac{4}{9} + \frac{4}{9} \cdot \frac{5}{9} = 1 - \left[ \left(\frac{5}{9}\right)^2 + \left(\frac{4}{9}\right)^2 \right]$$



# CART

→ 무조건 BInary split



## Binary Split

- Date = {Sun, Mon, Tue, Wed, Thu, Fri, Sat}
- Test all the split points
  - {Sun}, {Mon, Tue, Wed, Thu, Fri, Sat} ←
  - {Mon}, {Sun, Tue, Wed, Thu, Fri, Sat}
  - {Tue}, {Sun, Mon, Wed, Thu, Fri, Sat}
  - {Wed}, {Sun, Mon, Tue, Thu, Fri, Sat}
  - {Thu}, {Sun, Mon, Tue, Wed, Fri, Sat}
  - {Fri}, {Sun, Mon, Tue, Wed, Thu, Sat}
  - {Sat}, {Sun, Mon, Tue, Wed, Thu, Fri}

① Try all the split points

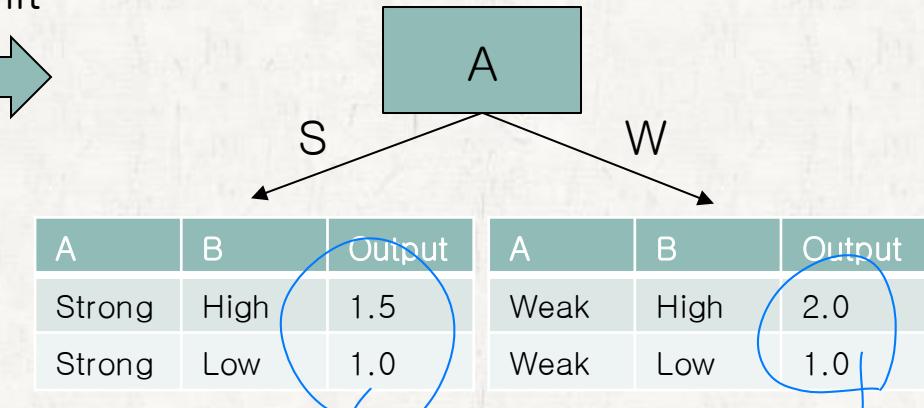
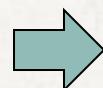
② choose the best "split point"

# Regression Tree

## How to Build a Regression Tree

A	B	Output
Strong	High	1.5
Strong	Low	1.0
Weak	High	2.0
Weak	Low	1.0

Split



Tree

1.25

1.25

1.5

The output of node n

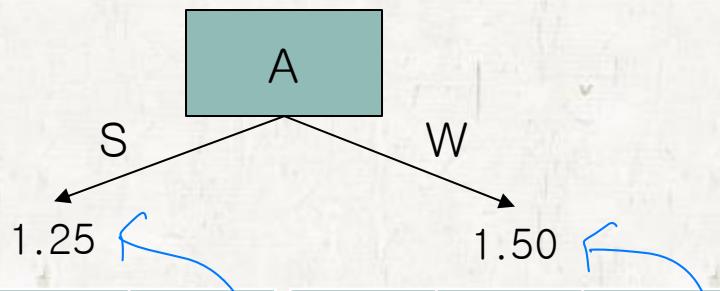
$$y_n = \frac{1}{|n|} \sum_{d \in n} t_d$$

# Regression Tree

1. Evaluate the **entropy** of the given table
  2. Split the table with an input
  3. Evaluate the **entropy** of each sub-table
  4. Evaluate the average of the **entropies**
  5. Evaluate the Gain: (1)–(4)
- \* Repeat this for all inputs and choose the best

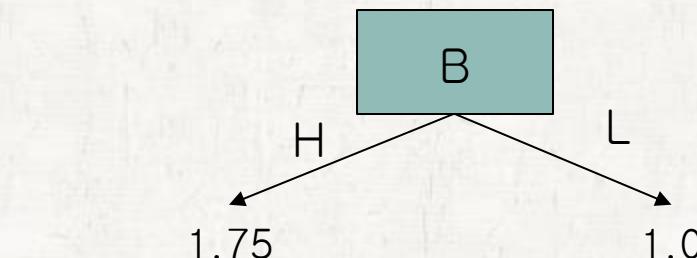
## How to Build a Tree

- Step 1: Split the dataset with respect to all the attributes



A	B	Output
Strong	High	1.5
Strong	Low	1.0

A	B	Output
Weak	High	2.0
Weak	Low	1.0



A	B	Output
Strong	High	1.5
Weak	High	2.0

A	B	Output
Strong	Low	1.0
Weak	Low	1.0

target: 1.5  
prediction: 1.25  $\rightarrow (1.5 - 1.25)^2$

t: 1.0  
prediction: 1.25  $\rightarrow (1.0 - 1.25)^2$

Which split is better?

Instead of Entropy, we use MSE

$$\Rightarrow \text{MSE} : \sum_{i=1}^n (E - Y)^2$$

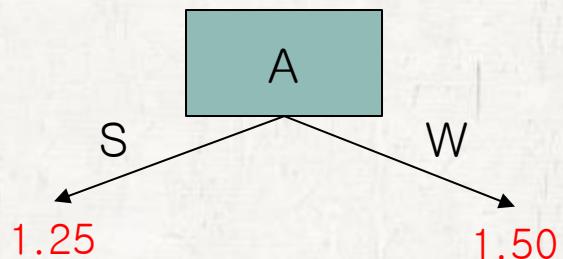
# Regression Tree

1. Split the table with an input
  2. Evaluate the **MSEs** of each sub-table
  3. Evaluate the sum of **MSEs**
- \* Repeat this for all inputs and choose the best

Entropy와 사용

## How to Build a Tree

- Step2: Evaluate the error of trees with respect to training data



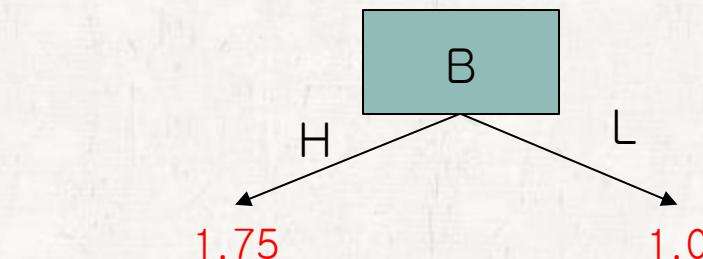
A	B	Output	A	B	Output
Strong	High	1.5	Weak	High	2.0
Strong	Low	1.0	Weak	Low	1.0

$$E = (1.5 - 1.25)^2 + (1.0 - 1.25)^2 \\ = 0.125$$

$$E = (2.0 - 1.5)^2 + (1.0 - 1.5)^2 \\ = 0.5$$

$$MSE = 0.125 + 0.5 = 0.625$$

둘을 비교



A	B	Output	A	B	Output
Strong	High	1.5	Strong	Low	1.0
Weak	High	2.0	Weak	Low	1.0

$$E = (1.5 - 1.75)^2 + (2.0 - 1.75)^2 \\ = 0.125$$

$$E = (1.0 - 1.0)^2 + (1.0 - 1.0)^2 \\ = 0.0$$

$$MSE = 0.125 + 0 = 0.125 \text{ less Error!}$$

$$MSE = \sum_{i \in L} (y_i - \bar{y}_L)^2 + \sum_{i \in R} (y_i - \bar{y}_R)^2$$

= 분할

# Regression Tree

## How to Build a Tree

- Step3: Choose the attr. with the smallest MSE, and repeat

A regression tree diagram where the root node is labeled 'B'. It splits into two branches: 'H' (left) and 'L' (right). Below the tree is a table with 6 data points:

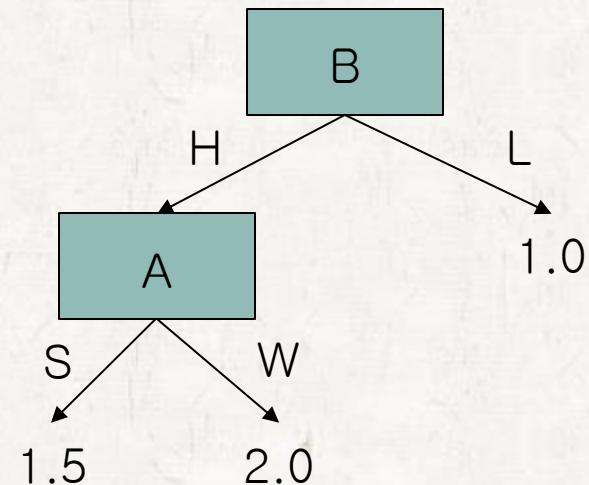
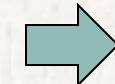
		A	B	Output
		Strong	High	1.5
		Weak	High	2.0
		Strong	Low	1.0
		Weak	Low	1.0

MSE가 0이 될 때까지 반복

The tree has split into node 'A' (left) and node 'B' (right). Node 'A' further splits into 'S' (left) and 'W' (right). Below the tree are two tables showing the resulting data splits:

		A	B	Out
		Strong	High	1.5
		Weak	High	2.0
		A	B	Out
		Strong	Low	1.0
		Weak	Low	1.0

Tree

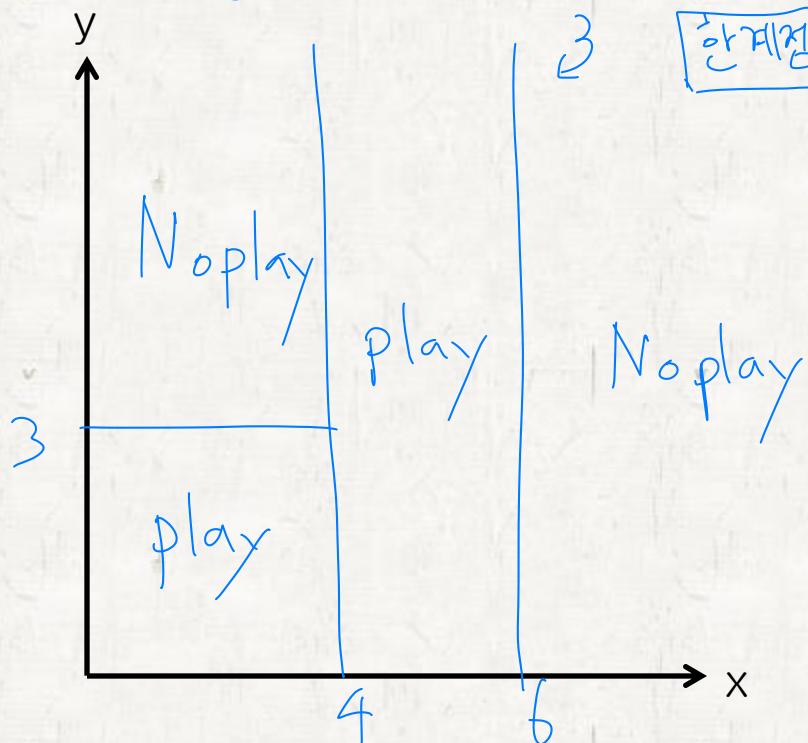




- How a Decision Splits Input Space

- Continuous attribute case

결정트리에 의한 입력 공간 분할 (수평, 수직 불가능)



실수형 attribute input

Input 1	Input 2	output
1	Y	play
10	5	No
5	7	No
6	10	yes
:	:	:

