

**기계학습 (2022년도 2학기)**

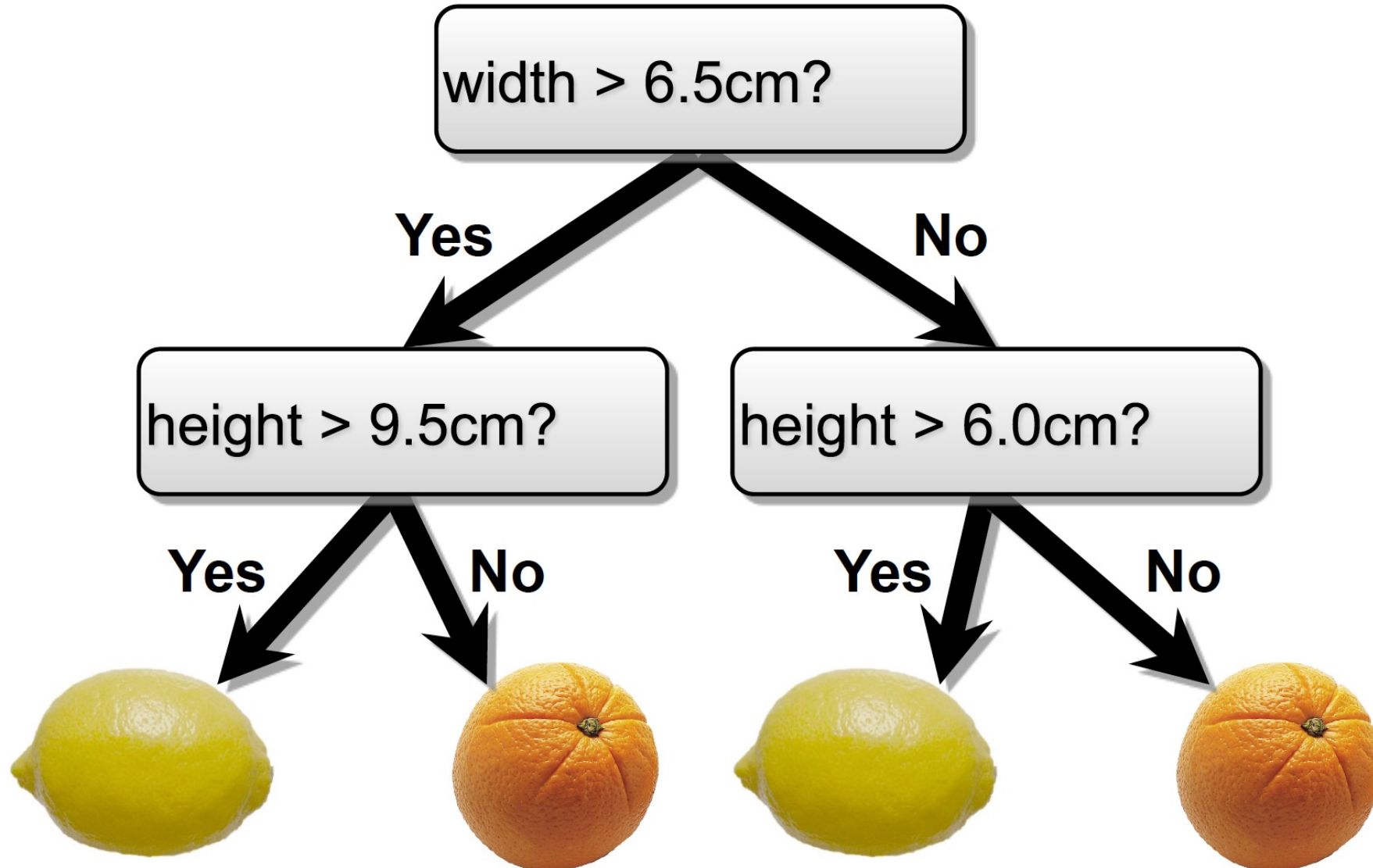
**Decision Trees**

**전북대학교 컴퓨터공학부**

# Today

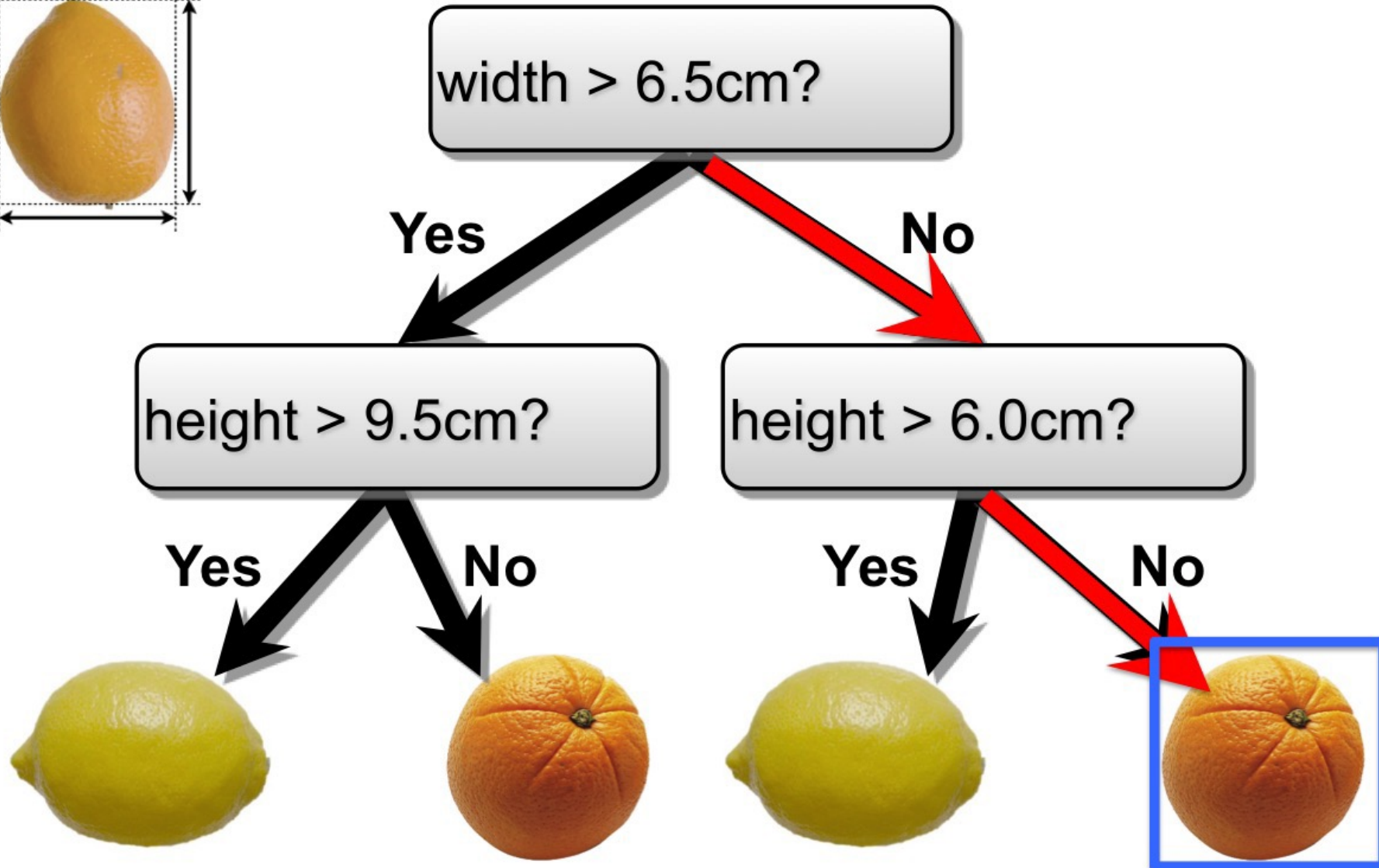
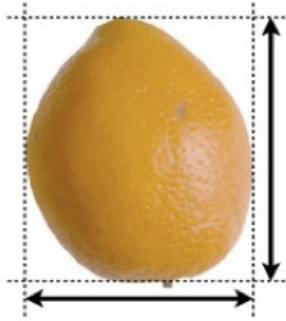
- Decision Trees
  - Simple but powerful learning algorithm
  - One of the most widely used learning algorithms in Kaggle competitions
- Useful information theoretic concepts (entropy, mutual information, etc.)

# Decision Trees



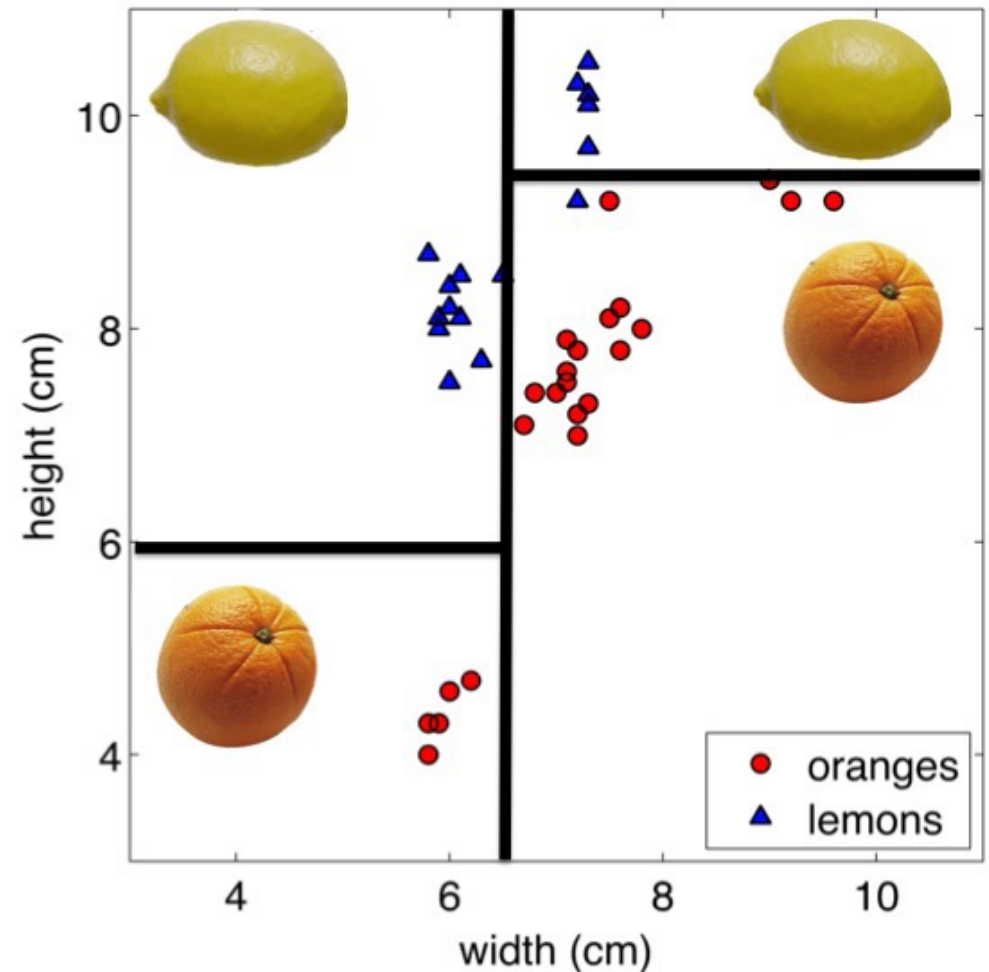
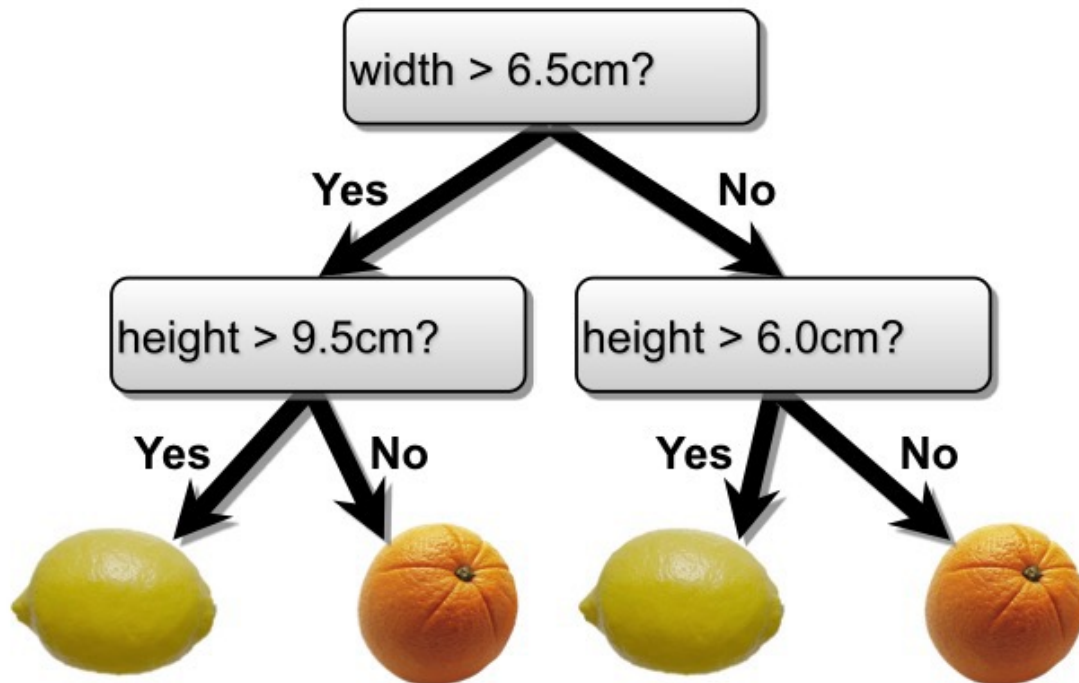
# Decision Trees

Test example



# Decision Trees

- Decision trees make predictions by recursively splitting on different attributes according to a tree structure.



# When to Consider Decision Trees

- Instances describable by attribute-value pairs
- Target function is discrete valued
- Disjunctive hypothesis may be required
- Possibly noisy training data

Examples:

- Equipment or medical diagnosis
- Credit risk analysis
- Modeling calendar scheduling preferences

# Example with Discrete Inputs

## ■ What if the attributes are discrete?

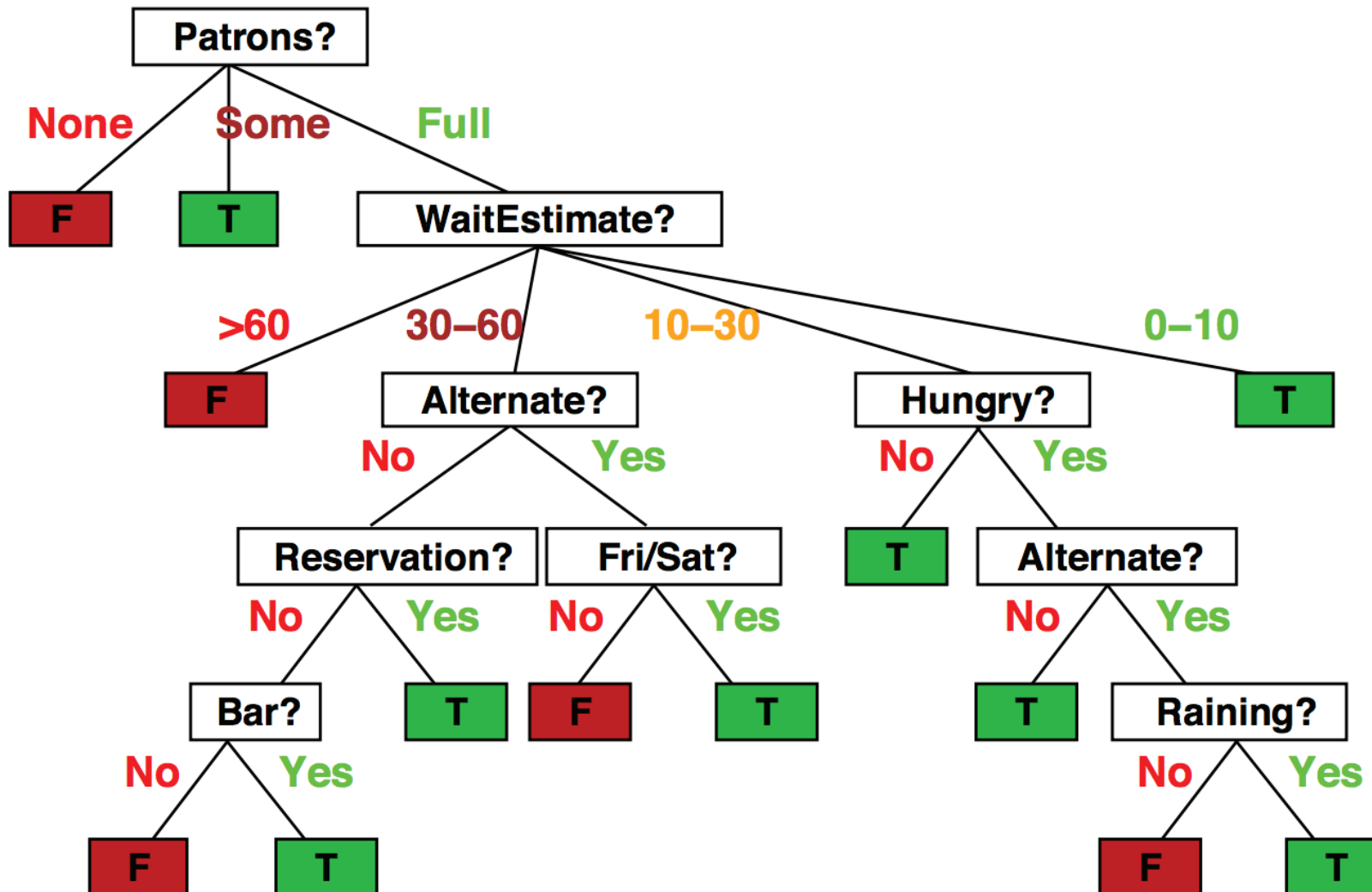
Example	Input Attributes										Goal
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
$x_1$	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0–10	$y_1 = \text{Yes}$
$x_2$	Yes	No	No	Yes	Full	\$	No	No	Thai	30–60	$y_2 = \text{No}$
$x_3$	No	Yes	No	No	Some	\$	No	No	Burger	0–10	$y_3 = \text{Yes}$
$x_4$	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10–30	$y_4 = \text{Yes}$
$x_5$	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	$y_5 = \text{No}$
$x_6$	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0–10	$y_6 = \text{Yes}$
$x_7$	No	Yes	No	No	None	\$	Yes	No	Burger	0–10	$y_7 = \text{No}$
$x_8$	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0–10	$y_8 = \text{Yes}$
$x_9$	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	$y_9 = \text{No}$
$x_{10}$	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10–30	$y_{10} = \text{No}$
$x_{11}$	No	No	No	No	None	\$	No	No	Thai	0–10	$y_{11} = \text{No}$
$x_{12}$	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30–60	$y_{12} = \text{Yes}$

*Attributes:*

1.	Alternate: whether there is a suitable alternative restaurant nearby.
2.	Bar: whether the restaurant has a comfortable bar area to wait in.
3.	Fri/Sat: true on Fridays and Saturdays.
4.	Hungry: whether we are hungry.
5.	Patrons: how many people are in the restaurant (values are None, Some, and Full).
6.	Price: the restaurant's price range (\$, \$\$, \$\$\$).
7.	Raining: whether it is raining outside.
8.	Reservation: whether we made a reservation.
9.	Type: the kind of restaurant (French, Italian, Thai or Burger).
10.	WaitEstimate: the wait estimated by the host (0-10 minutes, 10-30, 30-60, >60).

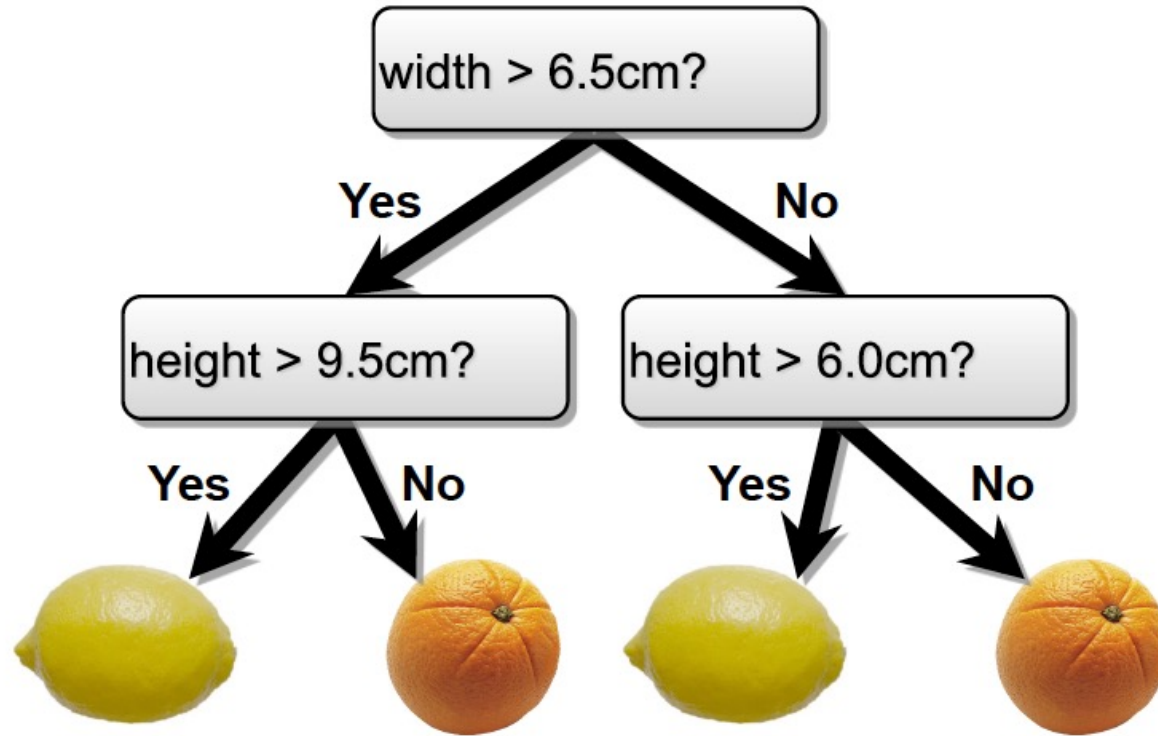
# Decision Tree: Example with Discrete Inputs

- The tree to decide whether to wait (T) or not (F)





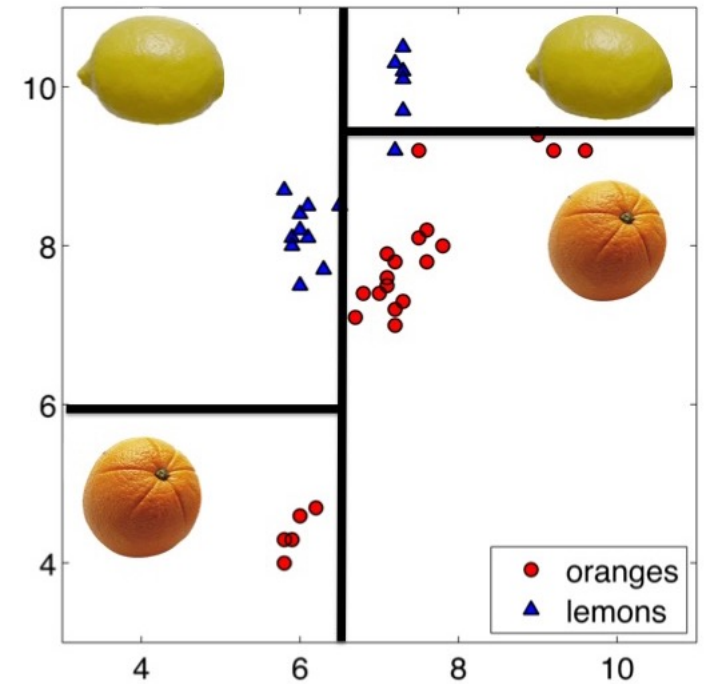
# Decision Trees



- Internal nodes test attributes
- Branching is determined by attribute value
- Leaf nodes are outputs (predictions)

# Decision Tree: Classification and Regression

- Each path from root to a leaf defines a region  $R_m$  of input space
- Let  $\{(x^{(m_1)}, t^{(m_1)}), \dots, (x^{(m_k)}, t^{(m_k)})\}$  be the training examples that fall into  $R_m$
- **Classification tree:**
  - discrete output
  - leaf value  $y_m$  typically set to the most common value in  $\{t^{(m_1)}, \dots, t^{(m_k)}\}$
- **Regression tree:**
  - continuous output
  - leaf value  $y_m$  typically set to the mean value in  $\{t^{(m_1)}, \dots, t^{(m_k)}\}$

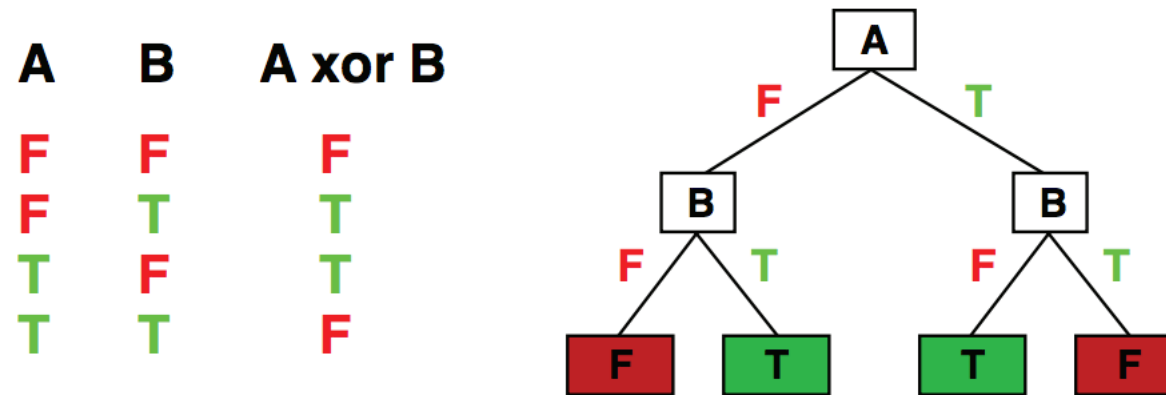


*Note: We will focus on classification*

# Expressiveness

## ■ Discrete-input, discrete-output case:

- Decision trees can express any function of the input attributes
- E.g., for Boolean functions, truth table row  $\rightarrow$  path to leaf:



## ■ Continuous-input, continuous-output case:

- Can approximate any function arbitrarily closely

## ■ Trivially, there is a consistent decision tree for any training set with one path to leaf for each example (unless $f$ nondeterministic in $x$ ) but it probably won't generalize to new examples

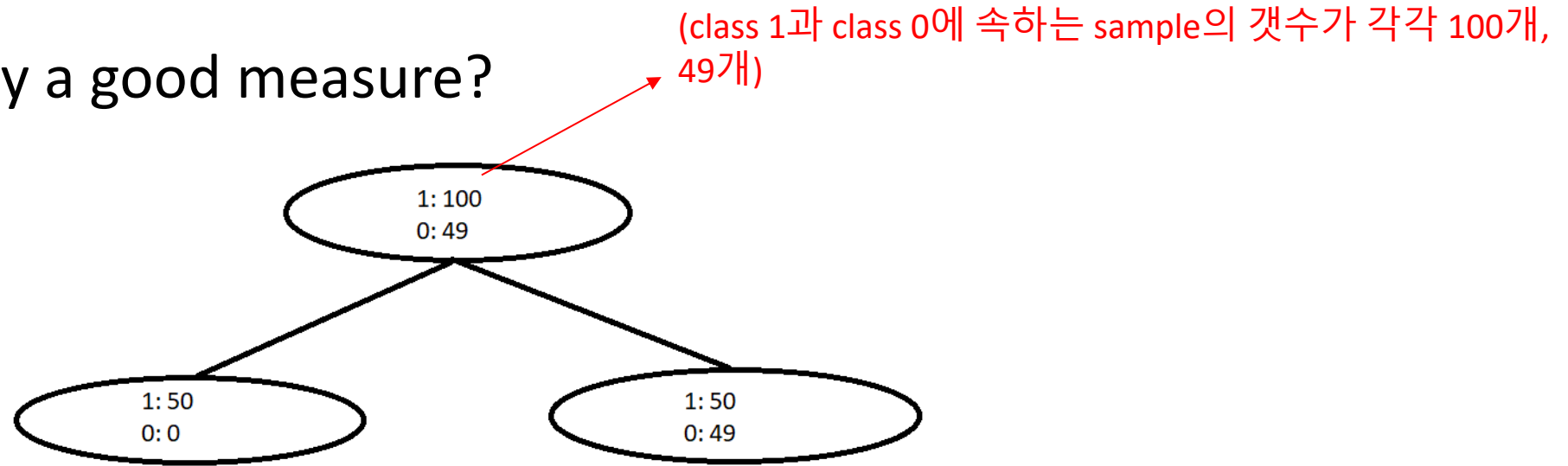
(모든 training data sample들이 별도의 leaf node를 가지도록 decision tree를 만들 수 있으나 general하지 않음)

# How do we Learn a Decision Tree?

- How do we construct a useful decision tree?
- Learning the simplest (smallest) decision tree is an NP complete problem [if you are interested, check: Hya I & Rivest'76]
  - Resort to a **greedy heuristic**:
    - Start from an empty decision tree
    - Split on the “best” attribute
    - Recurse
  - Which attribute is the “best”?
    - Choose based on accuracy?

# Choosing a Good Split

- Why isn't accuracy a good measure?



- Is this split good? Zero accuracy gain.  
(윗쪽 node만 이용한 분류와 아래쪽 node 2개를 사용한 분류 결과는 모두 1로 동일)
- Instead, we will use techniques from **information theory**

**Idea:** Use counts at leaves to define probability distributions, so we can measure uncertainty

# Choosing a Good Split

- Which attribute is better to split on,  $X_1$  or  $X_2$ ?
  - Deterministic: good (all are true or false; just one class in the leaf)
  - Uniform distribution: bad (all classes in leaf equally probable)
  - What about distributions in between?

Note: Let's take a slight detour and remember concepts from information theory

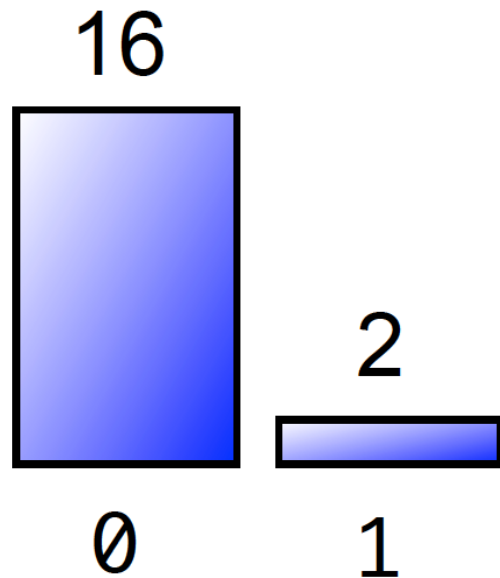
## We Flip Two Different Coins

- Sequence 1:

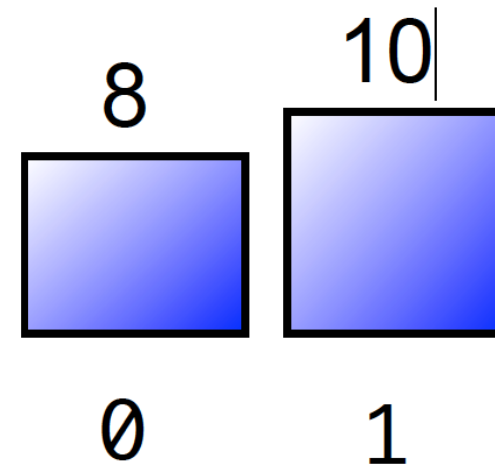
0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 ... ?

- Sequence 2:

0 1 0 1 0 1 1 1 0 1 0 0 1 1 0 1 0 1 ... ?



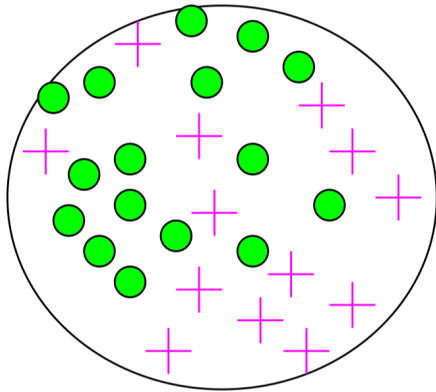
versus



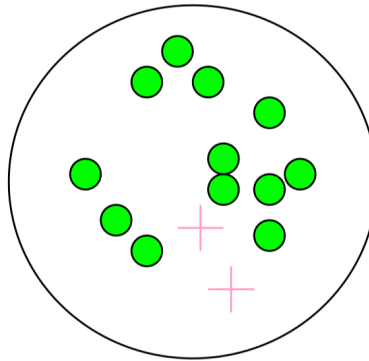
# Impurity/Entropy (informal)

- Measures the level of impurity in a group of examples

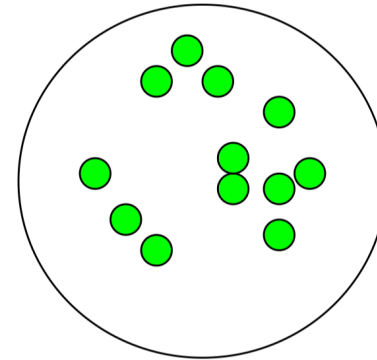
**Very impure group**



**Less impure**



**Minimum  
impurity**



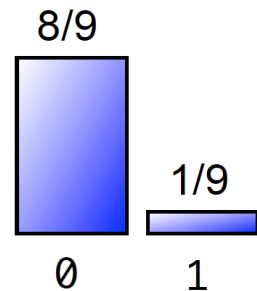


# Quantifying Uncertainty

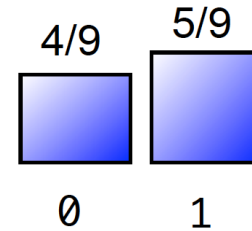
- **Entropy** is a measure of expected “surprise”:

(정보이론(Information theory)에서 나옴. Entropy가 클수록 또는 불확실성이 큰 이벤트  
일수록 많은 정보를 가지고 있음)

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$



$$-\frac{8}{9} \log_2 \frac{8}{9} - \frac{1}{9} \log_2 \frac{1}{9} \approx \frac{1}{2}$$



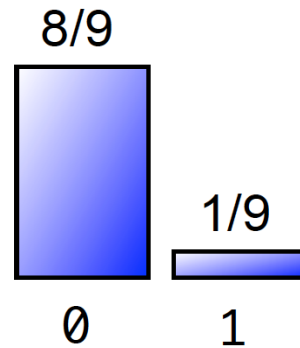
$$-\frac{4}{9} \log_2 \frac{4}{9} - \frac{5}{9} \log_2 \frac{5}{9} \approx 0.99$$

- Measures the information content of each observation
- Unit = **bits**
- A fair coin flip has 1 bit of entropy

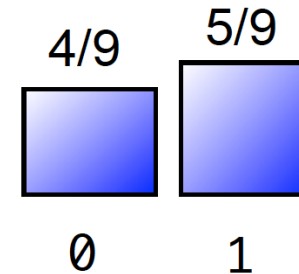
# Quantifying Uncertainty

- **Entropy** is a measure of expected “surprise”:  
(정보이론(Information theory)에서 나옴. Entropy가 클수록 또는 불확실성이 큰 이벤트 일수록 많은 정보를 가지고 있음)

$$H(X) = - \sum p(x) \log_2 p(x)$$



$$-\frac{8}{9} \log_2 \frac{8}{9} - \frac{1}{9} \log_2 \frac{1}{9} \approx \frac{1}{2}$$



$$-\frac{4}{9} \log_2 \frac{4}{9} - \frac{5}{9} \log_2 \frac{5}{9} \approx 0.99$$

- Measures the information content of each observation
- Unit = **bits**
- A fair coin flip has 1 bit of entropy (앞뒤 확률이 동일함: 가장 불확실성이 큼)

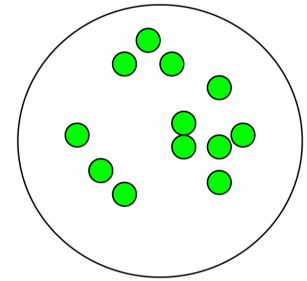
# What does Entropy mean for learning from examples?

## 2-Class Cases

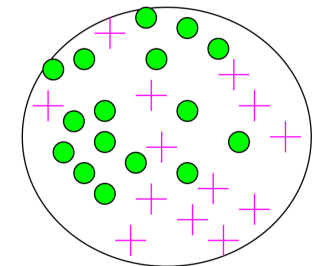
$$\text{Entropy } H(x) = - \sum_{i=1}^n P(x = i) \log_2 P(x = i)$$

- What is the entropy of a group in which all examples belong to the same class?
  - entropy =  $-(1 \times \log_2 1 + 0 \times \log_2 0) = 0$   
**not a good training set for learning**
- What is the entropy of a group with 50% in either class?
  - entropy =  $-(0.5 \times \log_2 0.5 + 0.5 \log_2 0.5) = 1$   
**good training set for learning**

**Minimum impurity**

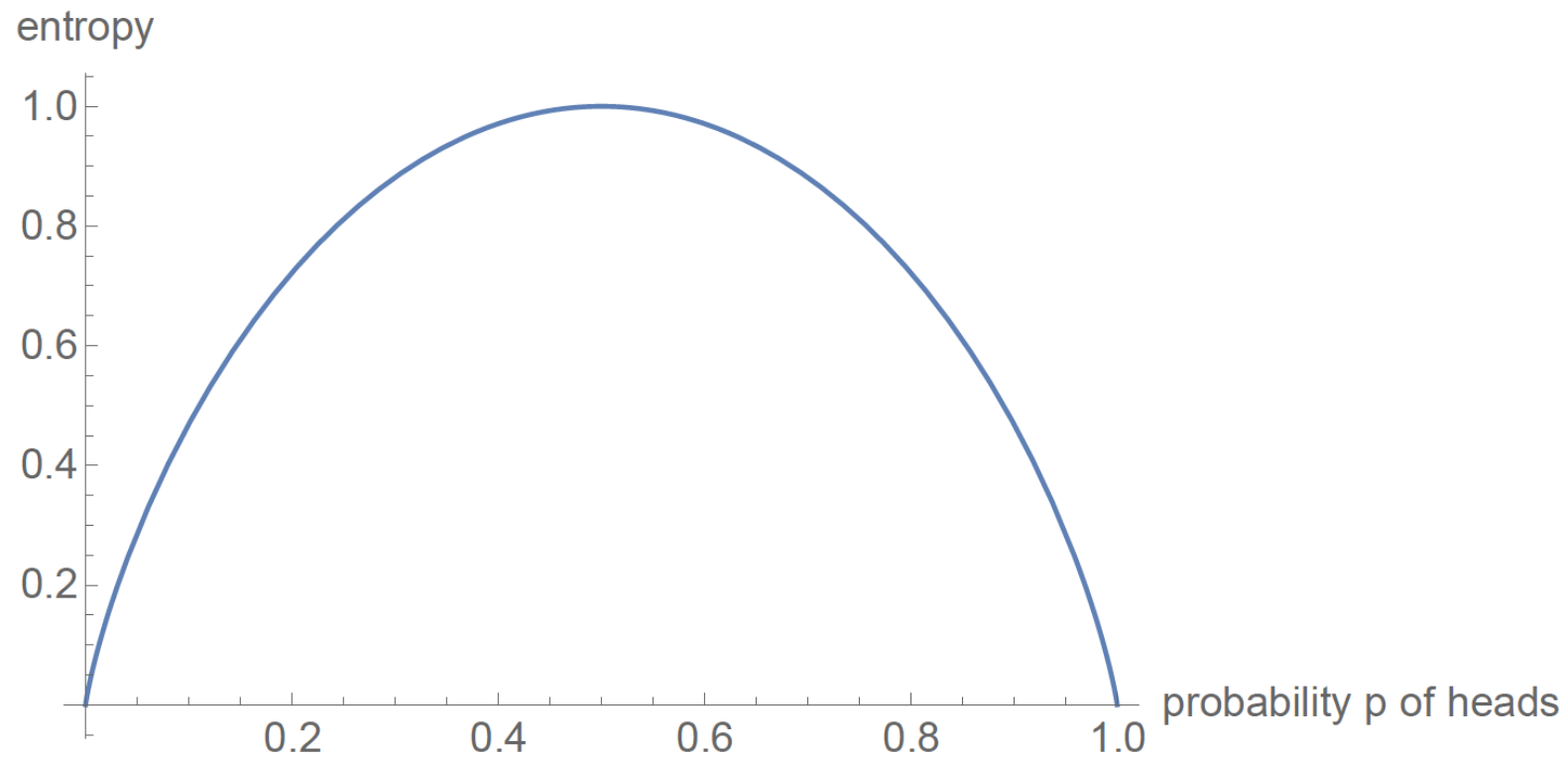


**Very impure group**



# Quantifying Uncertainty

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$



# Entropy

- “High Entropy”:
  - Variable has a uniform like distribution
  - Flat histogram
  - Values sampled from it are less predictable
- “Low Entropy”
  - Distribution of variable has many peaks and valleys
  - Histogram has many lows and highs
  - Values sampled from it are more predictable

# Entropy of a Joint Distribution

- Example:  $X = \{\text{Raining}, \text{Not raining}\}$ ,  $Y = \{\text{Cloudy}, \text{Not cloudy}\}$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

$$\begin{aligned} H(X, Y) &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y) \\ &= -\frac{24}{100} \log_2 \frac{24}{100} - \frac{1}{100} \log_2 \frac{1}{100} - \frac{25}{100} \log_2 \frac{25}{100} - \frac{50}{100} \log_2 \frac{50}{100} \\ &\approx 1.56 \text{bits} \end{aligned}$$

## Specific Conditional Entropy

- Example:  $X = \{\text{Raining}, \text{Not raining}\}$ ,  $Y = \{\text{Cloudy}, \text{Not cloudy}\}$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

- What is the entropy of cloudiness  $Y$ , given that it is raining?

$$\begin{aligned} H(Y|X = x) &= - \sum_{y \in Y} p(y|x) \log_2 p(y|x) \\ &= - \frac{24}{25} \log_2 \frac{24}{25} - \frac{1}{25} \log_2 \frac{1}{25} \\ &\approx 0.24 \text{bits} \end{aligned}$$

- We used:  $p(y|x) = \frac{p(x,y)}{p(x)}$ , and  $p(x) = \sum_y p(x,y)$  (sum in a row)

# Conditional Entropy

- Example:  $X = \{\text{Raining}, \text{Not raining}\}$ ,  $Y = \{\text{Cloudy}, \text{Not cloudy}\}$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

- The expected conditional entropy:

$$\begin{aligned} H(Y|X) &= \sum_{x \in X} p(x) H(Y|X = x) \\ &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(y|x) \end{aligned}$$



# Conditional Entropy

- Example:  $X = \{\text{Raining}, \text{Not raining}\}$ ,  $Y = \{\text{Cloudy}, \text{Not cloudy}\}$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

- What is the entropy of cloudiness, given the knowledge of whether or not it is raining?

$$\begin{aligned} H(Y|X) &= \sum_{x \in X} p(x) H(Y|X = x) \\ &= \frac{1}{4} H(\text{cloudy} | \text{is raining}) + \frac{3}{4} H(\text{cloudy} | \text{not raining}) \\ &\approx 0.75 \text{ bits} \end{aligned}$$

# Some useful properties of Conditional Entropy

- $H$  is always non-negative

- Chain rule:

$$H(X, Y) = H(X|Y) + H(Y) = H(Y|X) + H(X)$$

- If  $X$  and  $Y$  independent, then  $X$  doesn't tell us anything about  $Y$ :

$$H(Y|X) = H(Y)$$

- But  $X$  tells us everything about  $Y$

$$H(Y|X) = 0$$

- By knowing  $X$ , we can only decrease uncertainty about  $Y$

$$H(Y|X) \leq H(Y)$$

# Information Gain

- Example:  $X = \{\text{Raining, Not raining}\}$ ,  $Y = \{\text{Cloudy, Not cloudy}\}$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

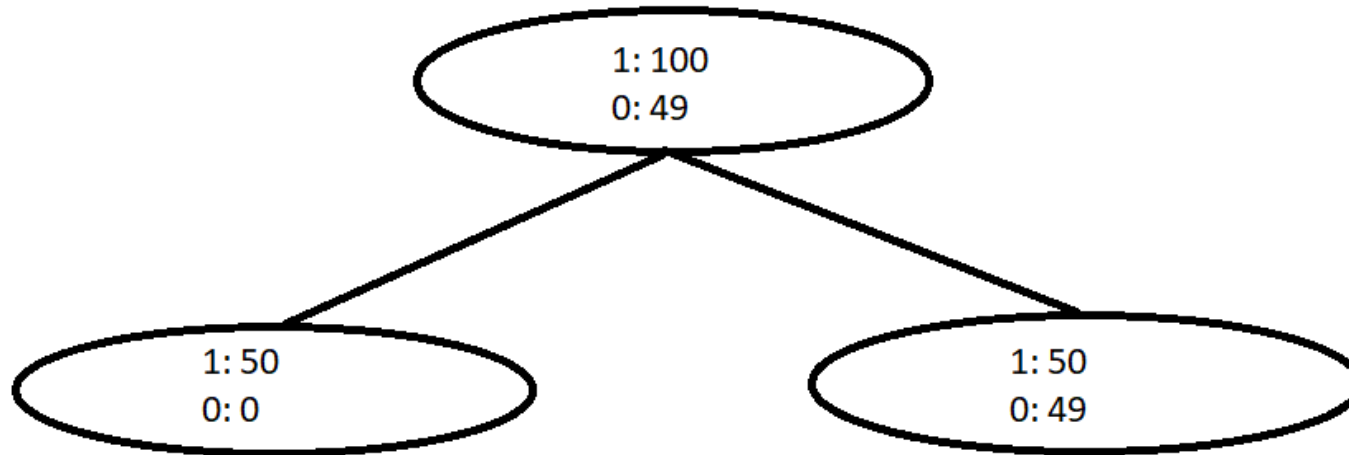
- How much information about cloudiness do we get by discovering whether it is raining?

$$\begin{aligned} IG(Y|X) &= H(Y) - H(Y|X) \\ &\approx 0.25 \text{ bits} \end{aligned}$$

- This is called the **information gain** in  $Y$  due to  $X$ , or the **mutual information** of  $Y$  and  $X$
- If  $X$  is completely uninformative about  $Y$ :  $IG(Y|X) = 0$
- If  $X$  is completely informative about  $Y$ :  $IG(Y|X) = H(Y)$

# Revisiting Our Original Example

- **Information Gain** =  $\text{entropy}(\text{parent}) - [\text{average entropy}(\text{children})]$
- Information gain measures the informativeness of a variable, which is exactly what we desire in a decision tree attribute!
- What is the information gain of this split?

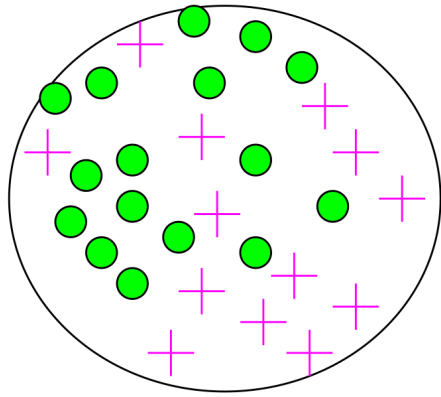


- Root entropy:  $H(Y) = -\frac{49}{149} \log_2\left(\frac{49}{149}\right) - \frac{100}{149} \log_2\left(\frac{100}{149}\right) \approx 0.91$
- Leafs entropy:  $H(Y|left) = 0$ ,  $H(Y|right) \approx 1$
- $IG(split) \approx 0.91 - \left(\frac{1}{3} \cdot 0 + \frac{2}{3} \cdot 1\right) \approx 0.24 > 0$

# More Example for Calculating Information Gain

- **Information Gain** = entropy(parent) – [average entropy(children)]

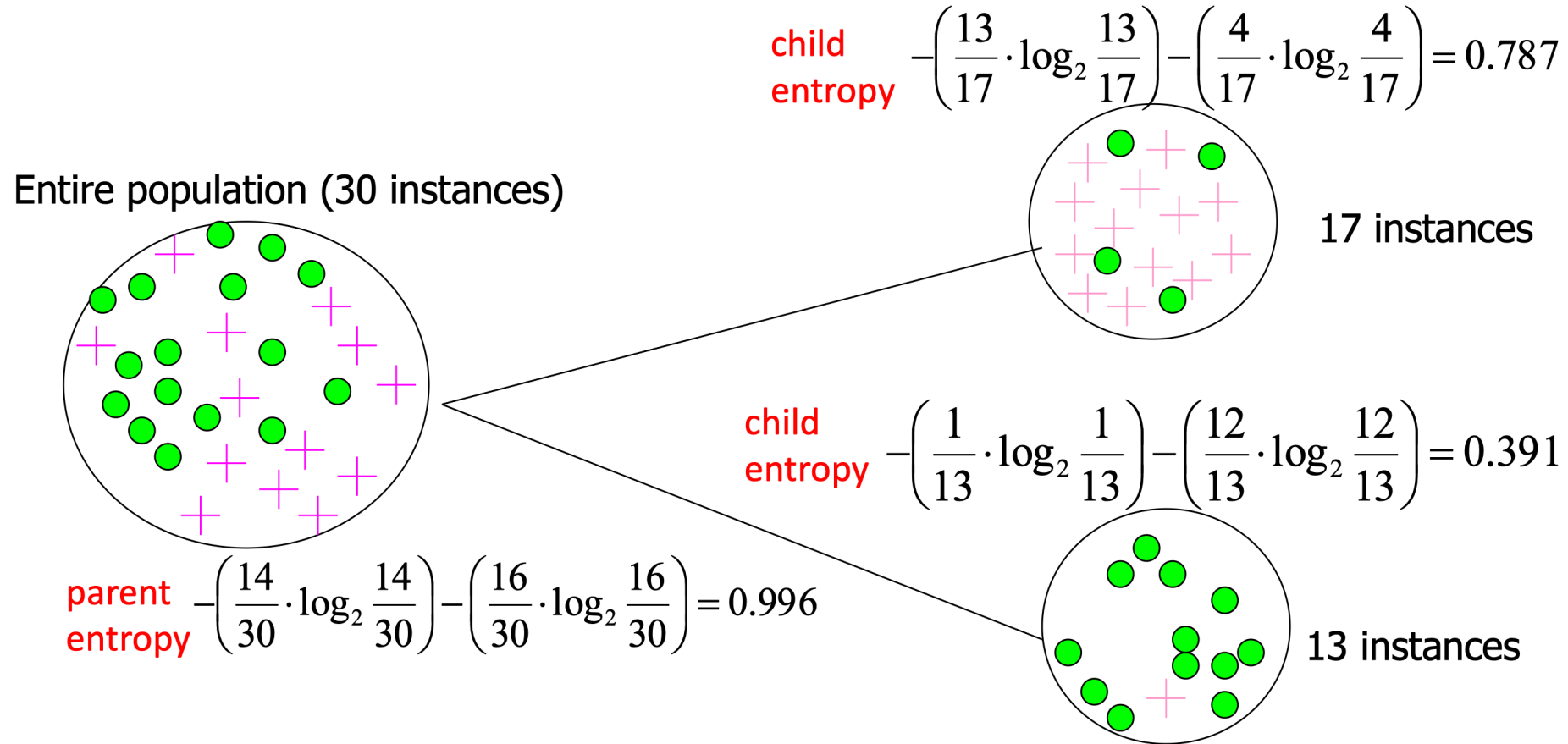
Entire population (30 instances)



$$\text{parent entropy} = -\left(\frac{14}{30} \cdot \log_2 \frac{14}{30}\right) - \left(\frac{16}{30} \cdot \log_2 \frac{16}{30}\right) = 0.996$$

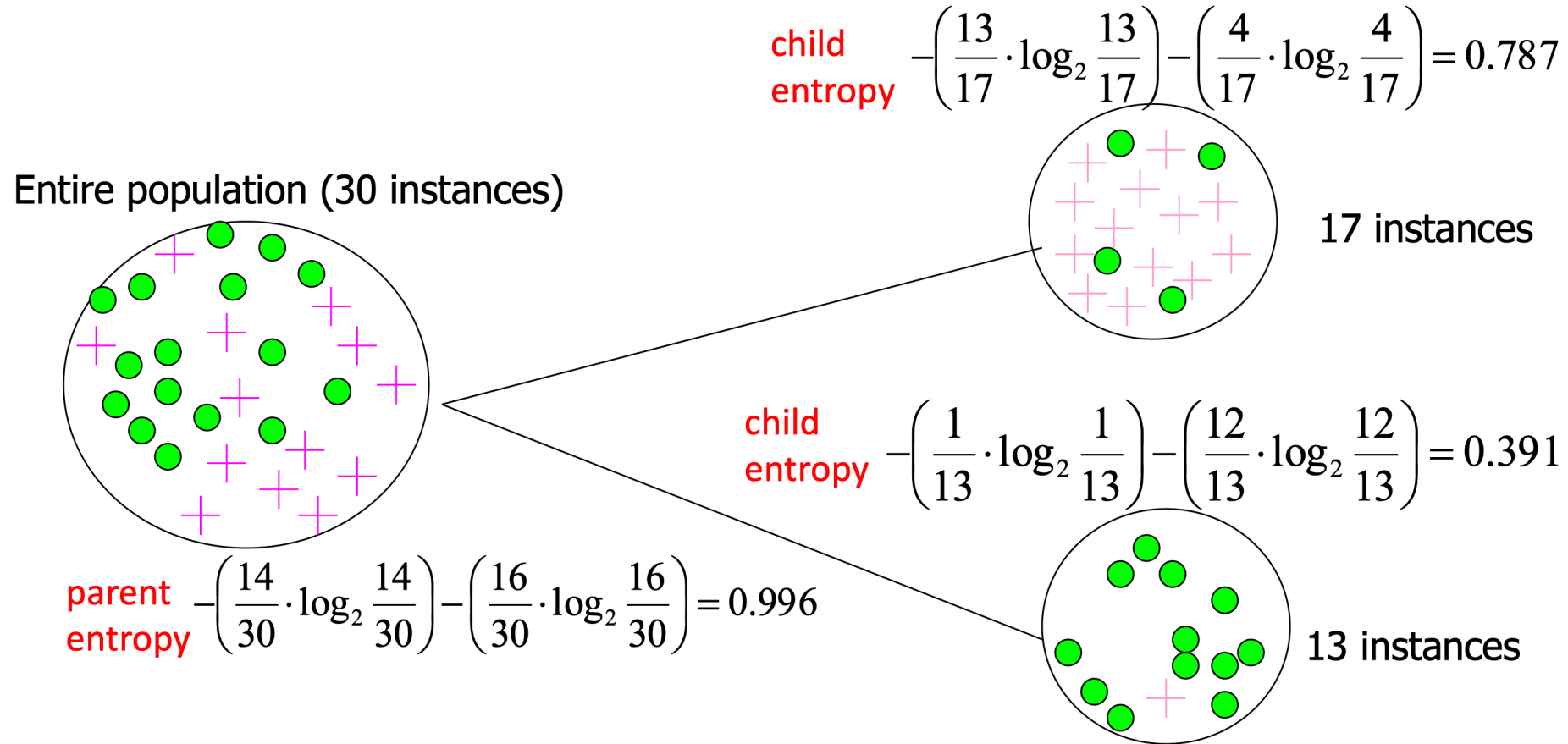
# More Example for Calculating Information Gain

- **Information Gain** = entropy(parent) – [average entropy(children)]



# More Example for Calculating Information Gain

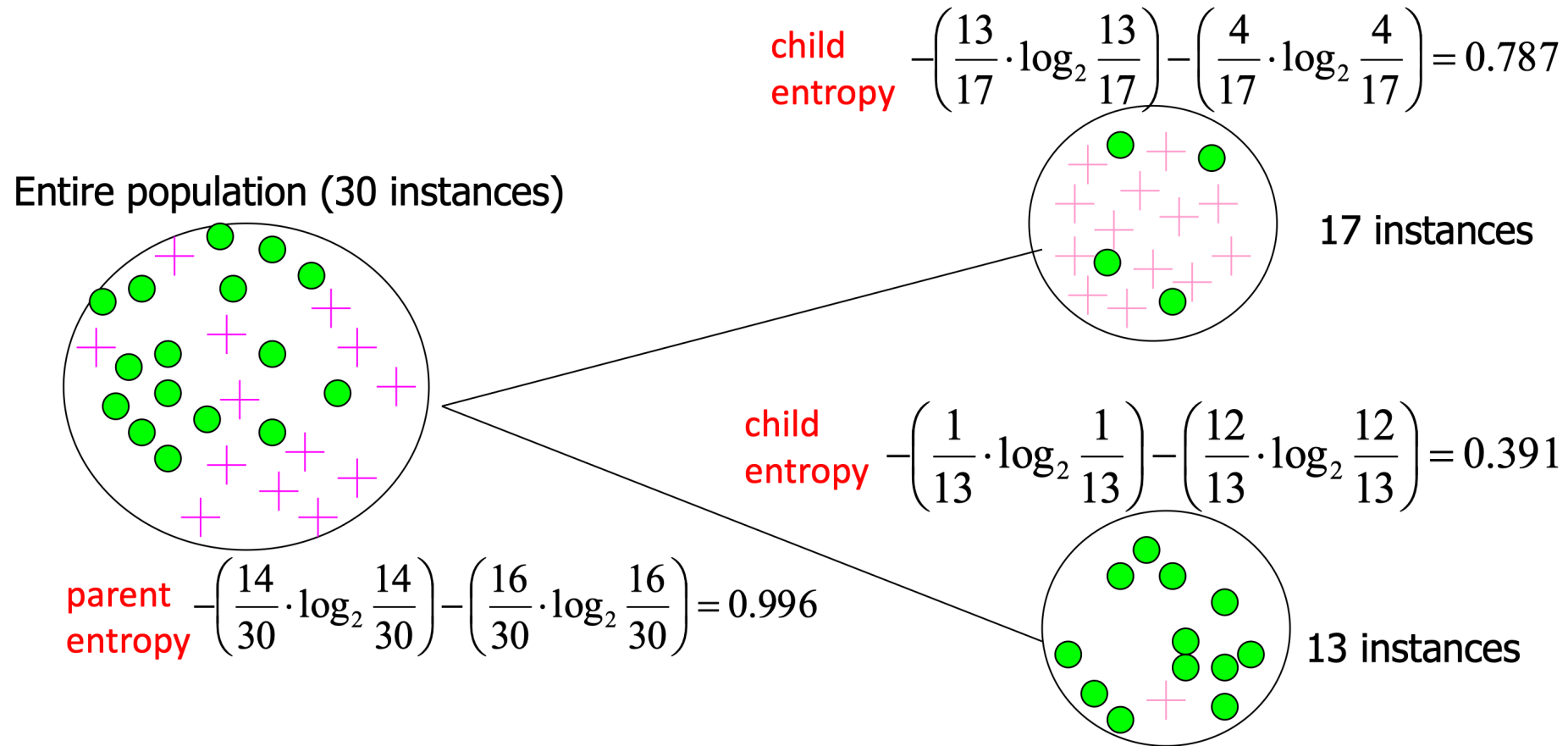
- **Information Gain** = entropy(parent) – [average entropy(children)]



- (Weighted) Average Entropy of Children =  $\left(\frac{17}{30} \times 0.787\right) + \left(\frac{13}{30} \times 0.391\right) = 0.615$

# More Example for Calculating Information Gain

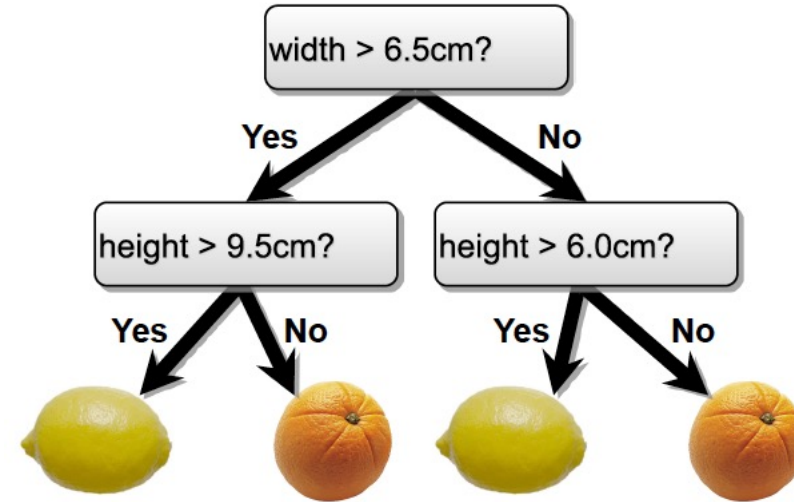
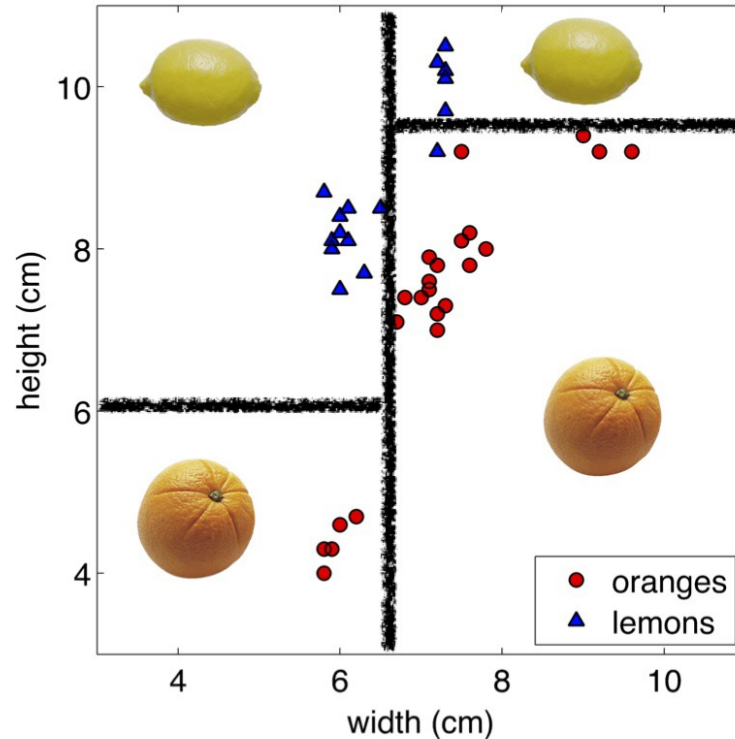
- **Information Gain** = entropy(parent) – [average entropy(children)]



- (Weighted) Average Entropy of Children =  $\left(\frac{17}{30} \times 0.787\right) + \left(\frac{13}{30} \times 0.391\right) = 0.615$
- Information gain =  $0.996 - 0.615 = 0.38$



# Constructing Decision Trees



- At each level, one must choose:
  1. Which variable to split.
  2. Possibly where to split it.
- Choose them based on how much information we would gain from the decision! (choose attribute that gives the best gain)

# Decision Tree Construction Algorithm

- Simple, greedy, recursive approach, builds up tree node-by-node
- 
1. pick an attribute to split at a non-terminal node
  2. split examples into groups based on attribute value
  3. for each group:
    - if no examples - return majority from parent
    - else if all examples in same class - return class
    - else loop to step 1

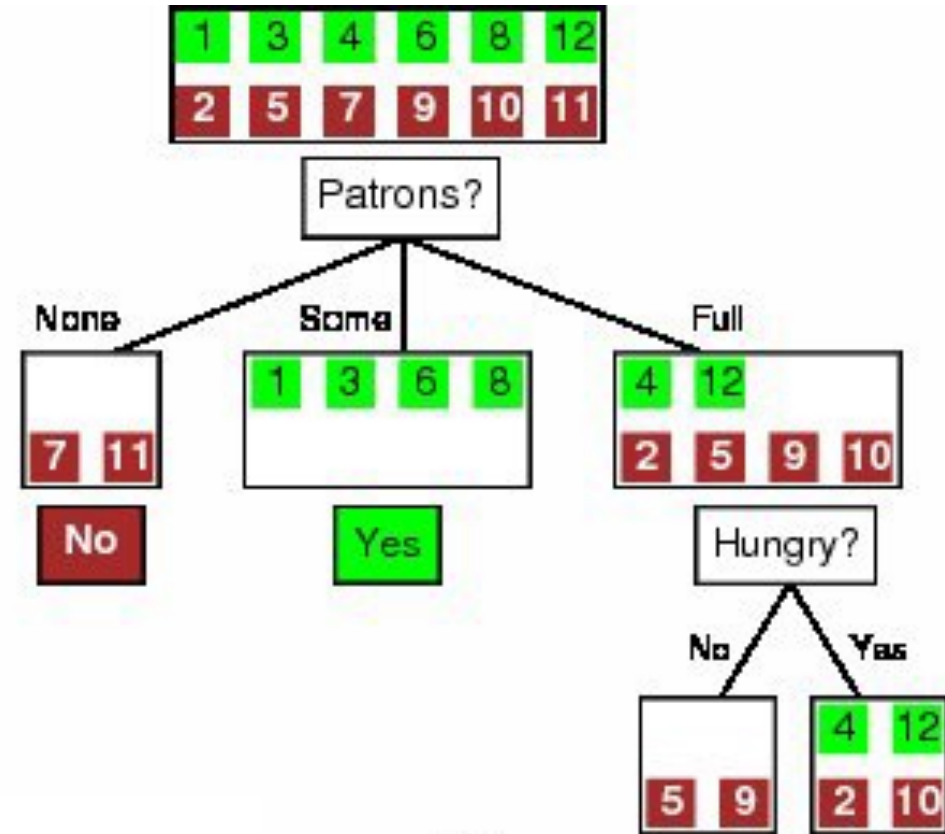
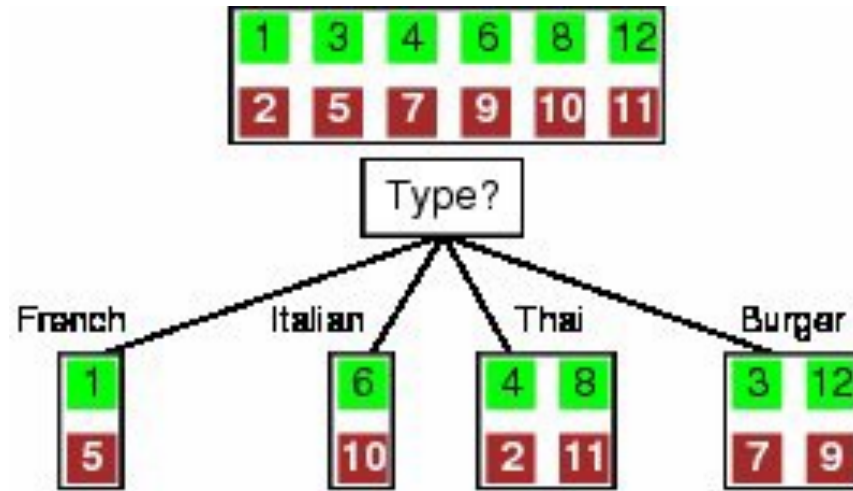
# Back to Our Example

Example	Input Attributes										Goal
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
$\mathbf{x}_1$	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0–10	$y_1 = \text{Yes}$
$\mathbf{x}_2$	Yes	No	No	Yes	Full	\$	No	No	Thai	30–60	$y_2 = \text{No}$
$\mathbf{x}_3$	No	Yes	No	No	Some	\$	No	No	Burger	0–10	$y_3 = \text{Yes}$
$\mathbf{x}_4$	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10–30	$y_4 = \text{Yes}$
$\mathbf{x}_5$	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	$y_5 = \text{No}$
$\mathbf{x}_6$	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0–10	$y_6 = \text{Yes}$
$\mathbf{x}_7$	No	Yes	No	No	None	\$	Yes	No	Burger	0–10	$y_7 = \text{No}$
$\mathbf{x}_8$	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0–10	$y_8 = \text{Yes}$
$\mathbf{x}_9$	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	$y_9 = \text{No}$
$\mathbf{x}_{10}$	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10–30	$y_{10} = \text{No}$
$\mathbf{x}_{11}$	No	No	No	No	None	\$	No	No	Thai	0–10	$y_{11} = \text{No}$
$\mathbf{x}_{12}$	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30–60	$y_{12} = \text{Yes}$

*Attributes:*

1.	Alternate: whether there is a suitable alternative restaurant nearby.
2.	Bar: whether the restaurant has a comfortable bar area to wait in.
3.	Fri/Sat: true on Fridays and Saturdays.
4.	Hungry: whether we are hungry.
5.	Patrons: how many people are in the restaurant (values are None, Some, and Full).
6.	Price: the restaurant's price range (\$, \$\$, \$\$\$).
7.	Raining: whether it is raining outside.
8.	Reservation: whether we made a reservation.
9.	Type: the kind of restaurant (French, Italian, Thai or Burger).
10.	WaitEstimate: the wait estimated by the host (0-10 minutes, 10-30, 30-60, >60).

# Attribute Selection

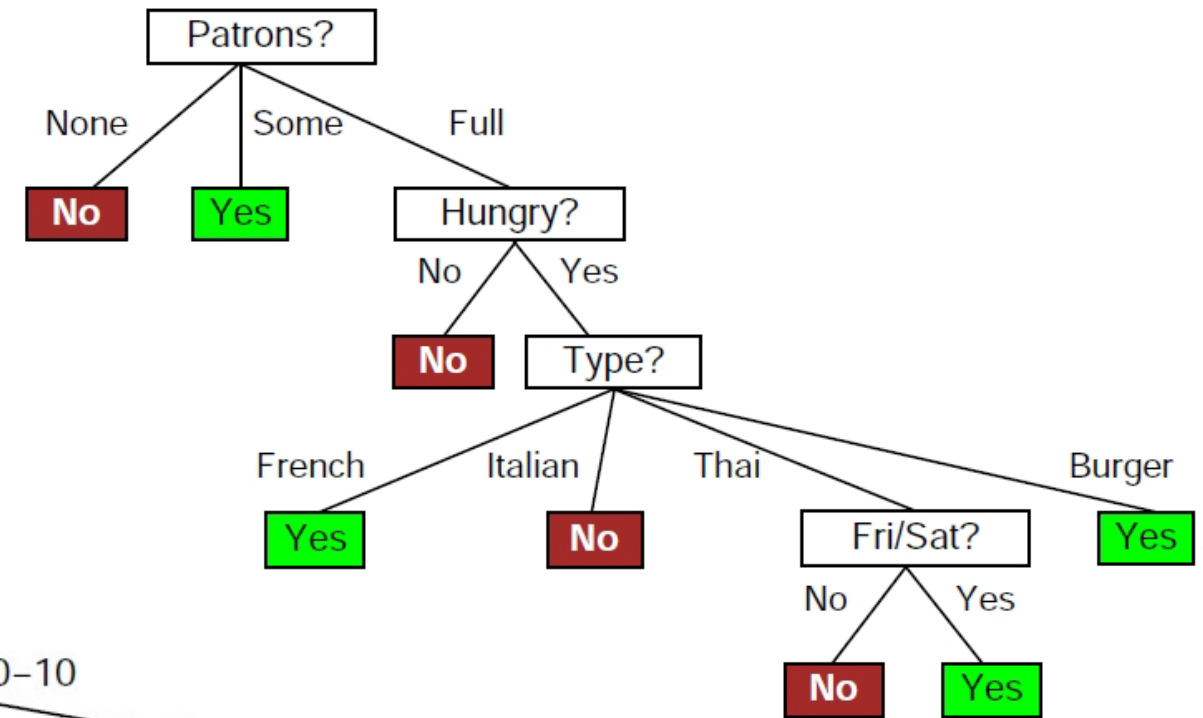
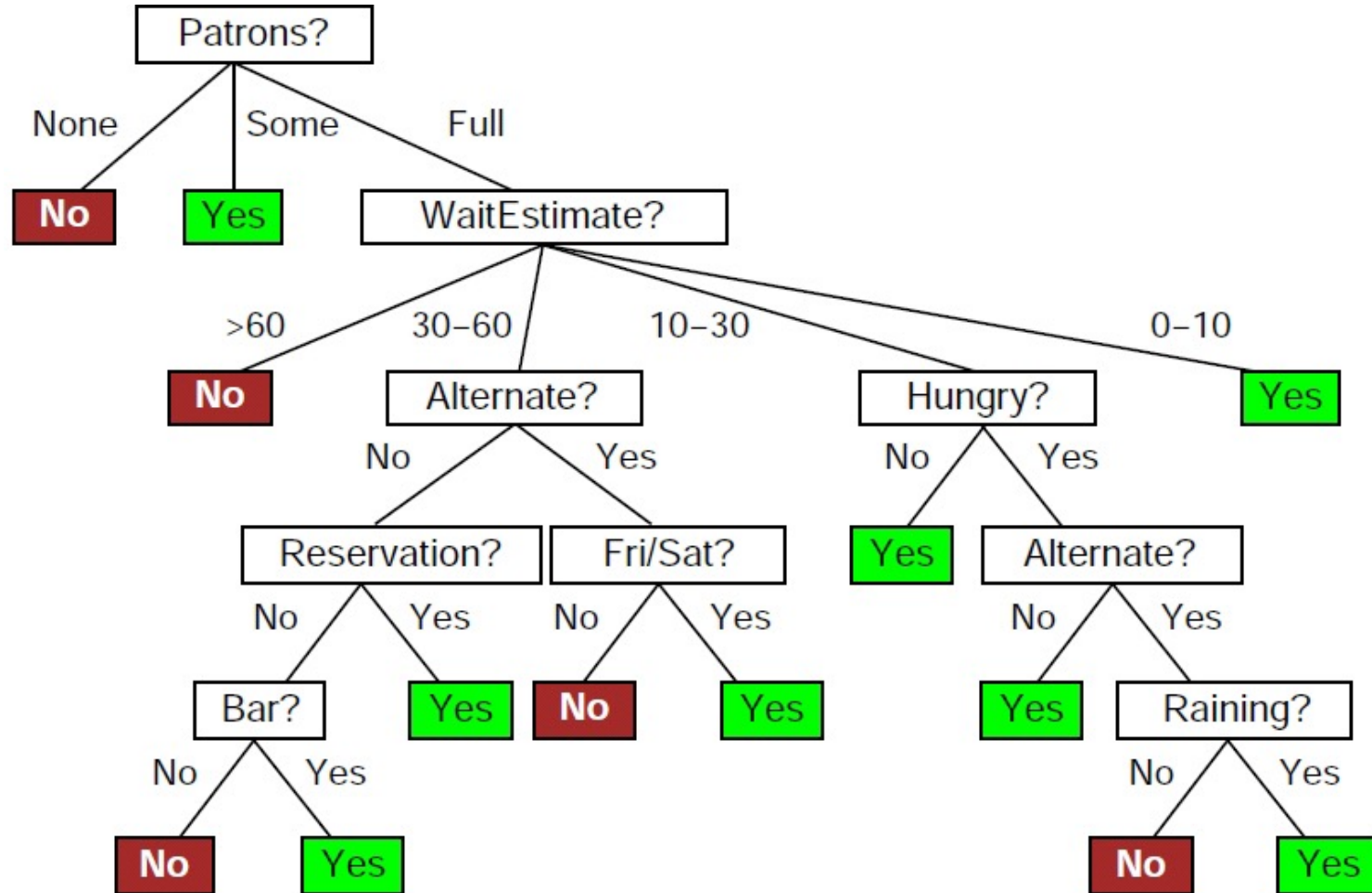


$$IG(Y) = H(Y) - H(Y|X)$$

$$IG(type) = 1 - \left[ \frac{2}{12} H(Y|Fr.) + \frac{2}{12} H(Y|It.) + \frac{4}{12} H(Y|Thai) + \frac{4}{12} H(Y|Bur.) \right] = 0$$

$$IG(Patrons) = 1 - \left[ \frac{2}{12} H(0, 1) + \frac{4}{12} H(1, 0) + \frac{6}{12} H\left(\frac{2}{6}, \frac{4}{6}\right) \right] \approx 0.541$$

# Which Tree is Better?



# What Makes a Good Tree?

- Not too small: need to handle important but possibly subtle distinctions in data
- Not too big:
  - Computational efficiency (avoid redundant, spurious attributes)
  - Avoid over-fitting training examples
  - Human interpretability
- **“Occam's Razor”**: find the simplest hypothesis that fits the observations
  - Useful principle, but hard to formalize (how to define simplicity?)
  - See Domingos, 1999, “The role of Occam's razor in knowledge discovery”
- We desire small trees with informative nodes near the root

# Decision Tree Miscellany

## ■ Problems:

- You have exponentially less data at lower levels
- Too big of a tree can overfit the data
- Greedy algorithms don't necessarily yield the global optimum

## ■ Handling continuous attributes

- Split based on a threshold, chosen to maximize information gain

## ■ Decision trees can also be used for regression on real-valued outputs. Choose splits to minimize squared error, rather than maximize information gain.

# Comparison to k-NN

## ■ Advantages of decision trees over k-NN

- Good with discrete attributes
- Easily deals with missing values (just treat as another value)  
(예: nominal feature, 즉 categorical feature인 경우 새로운 class로 취급)
- Robust to scale of inputs
- Fast at test time
- More interpretable

## ■ Advantages of k-NN over decision trees

- Able to handle attributes/features that interact in complex ways (e.g. pixels)
- Can incorporate interesting distance measures (e.g. shape contexts)
- Typically make better predictions in practice
  - As we'll see next lecture, ensembles of decision trees are much stronger. But they lose many of the advantages listed above.