

Decision Trees

Data Intelligence and Learning ([DIAL](#)) Lab

Prof. Jongwuk Lee

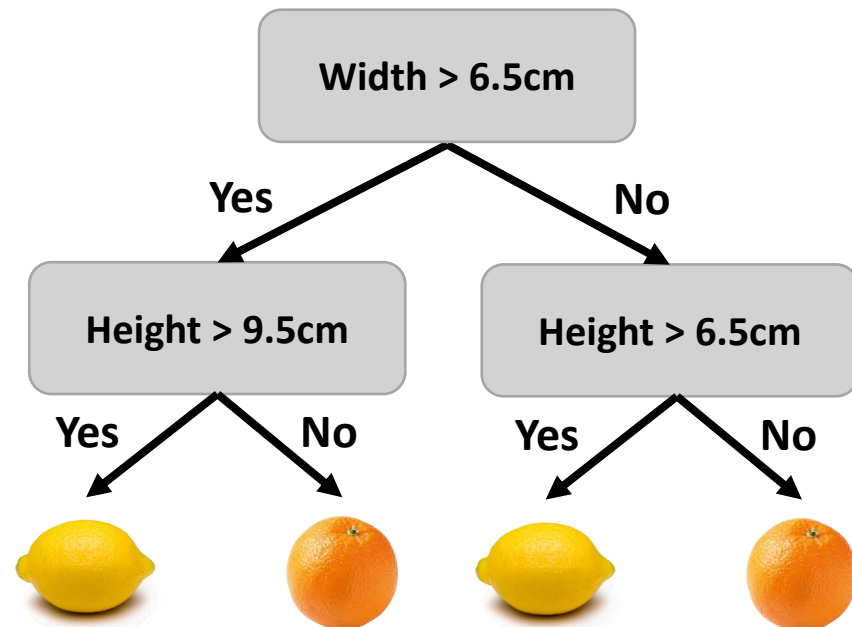


Decision Tree Basics

What are Decision Trees?



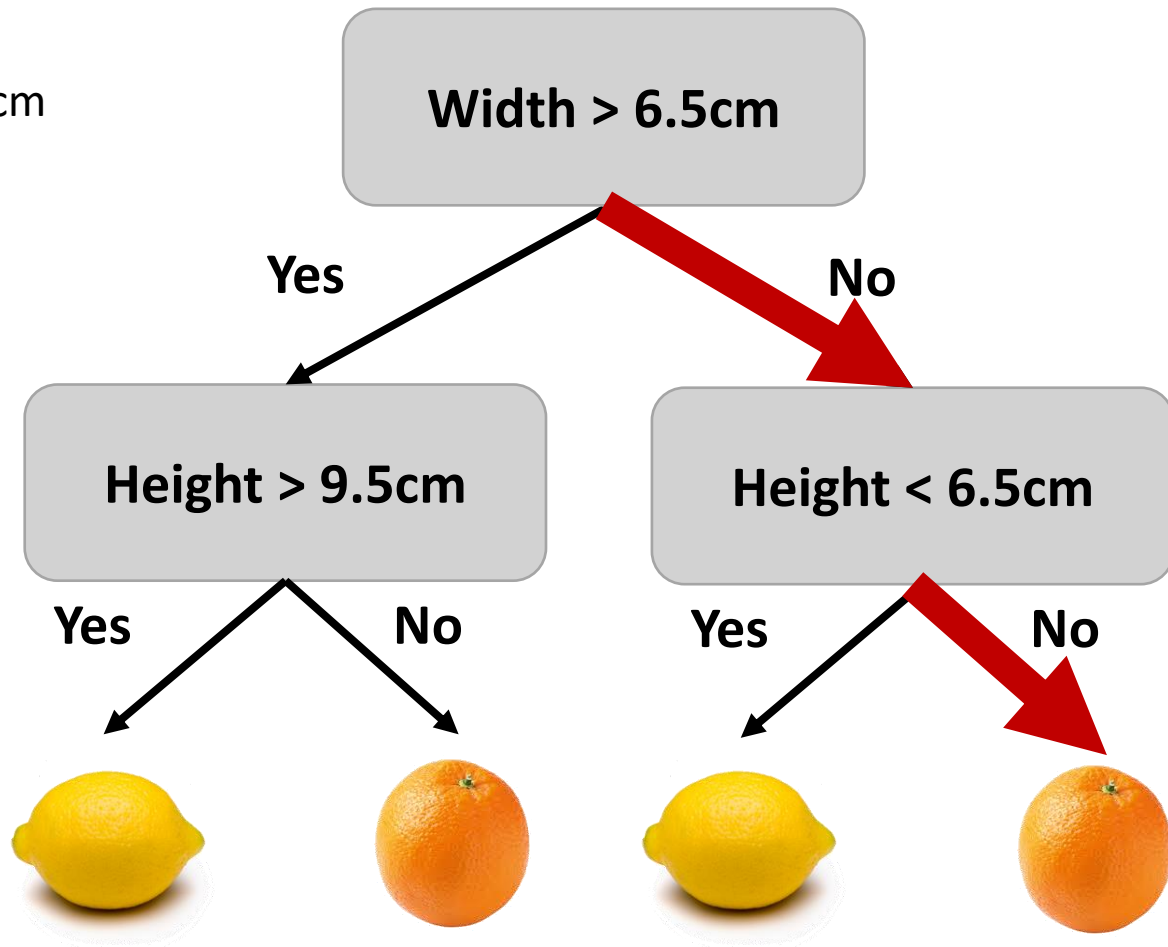
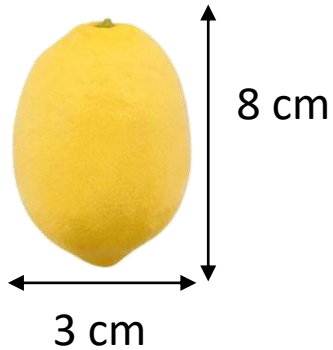
- Decision trees make predictions by recursively splitting into different attributes.
- ◆ Each **internal node** represents an **attribute** at each stage.
 - ◆ Each **branch** represents an **attribute value**.
 - ◆ Each **leaf node** represents a **class** label.
 - ◆ The path from the root to a leaf node represents a **classification rule**.



What are Decision Trees?



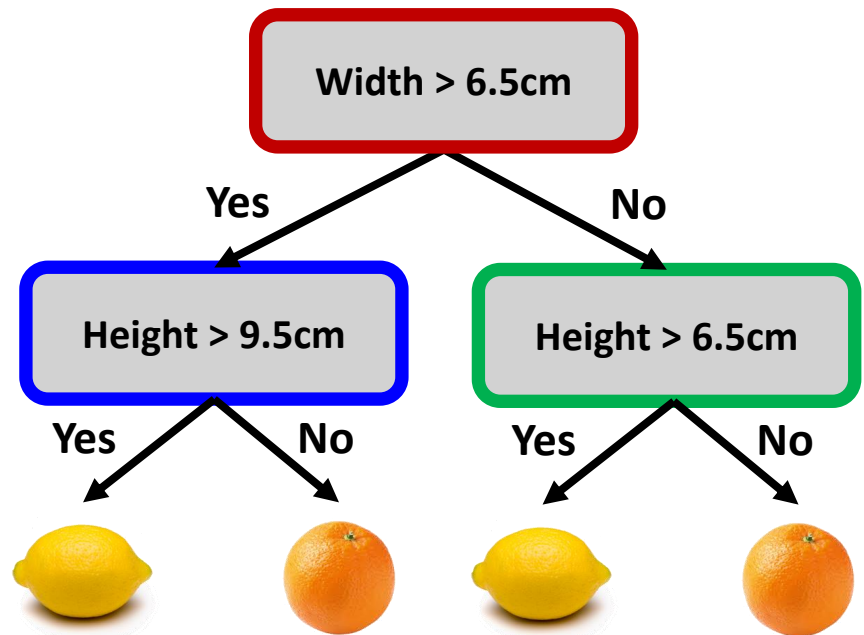
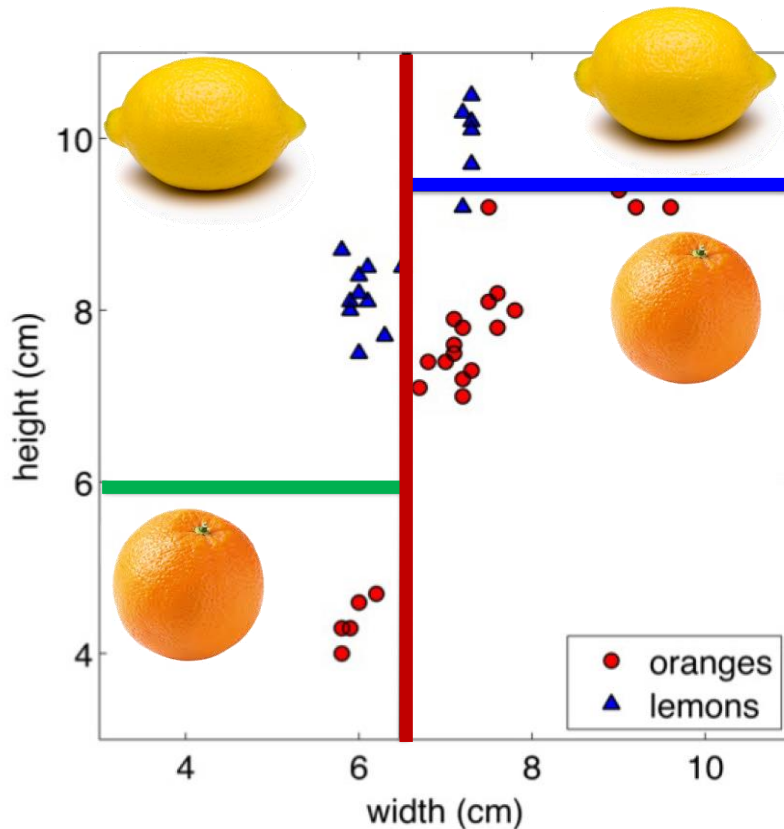
Test sample



What are Decision Trees?



- For a continuous attribute, we split it with a specific value.
 - ◆ The input space is **recursively** divided into **two regions** parallel to axes.



Classification and Regression



➤ Each path from the root to a leaf node corresponds to a **region R_m** of input space.

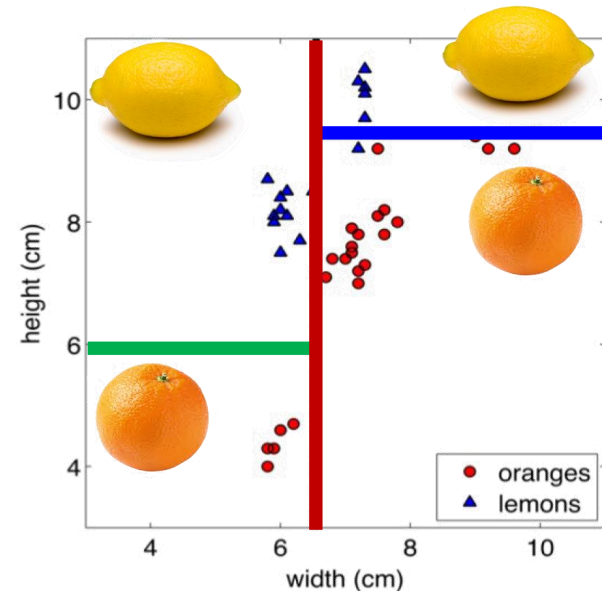
- ◆ Let $\{(x^{(m_1)}, y^{(m_1)}), \dots, (x^{(m_k)}, y^{(m_k)})\}$ be training samples in R_m .

➤ **Classification tree**

- ◆ Discrete output
- ◆ Leaf value represents the **most common value** in $\{y^{(m_1)}, \dots, y^{(m_k)}\}$.

➤ **Regression tree**

- ◆ Continuous output
- ◆ Leaf value represents the **mean value** in $\{y^{(m_1)}, \dots, y^{(m_k)}\}$.



Classification: Good vs. Bad



- We want to construct a **classification tree** to classify heroes as **good** or **bad** according to their appearance.



Training data

	<i>Gender</i>	<i>Mask</i>	<i>Cape</i>	<i>Tie</i>	<i>Ears</i>	<i>Smokes</i>	<i>Label</i>
Batman	Male	Yes	Yes	No	Yes	No	Good
Robin	Male	Yes	Yes	No	No	No	Good
Alfred	Male	No	No	Yes	No	No	Good
Penguin	Male	No	No	Yes	No	Yes	Bad
Catwoman	Female	Yes	No	No	Yes	No	Bad
Joker	Male	No	No	No	No	No	Bad



Classification: Good vs. Bad

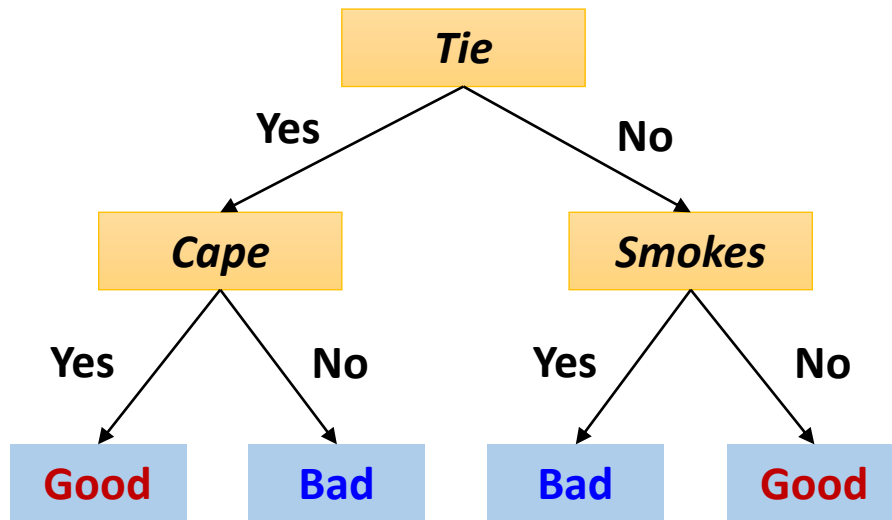
- A hero is a male who has a mask, a cape, and no tie.
- Q: Is he a good or bad man?
- How do we predict whether he is **good** or **bad**?

Training data

	<i>Gender</i>	<i>Mask</i>	<i>Cape</i>	<i>Tie</i>	<i>Ears</i>	<i>Smokes</i>	<i>Label</i>
Batman	Male	Yes	Yes	No	Yes	No	Good
Robin	Male	Yes	Yes	No	No	No	Good
Alfred	Male	No	No	Yes	No	No	Good
Penguin	Male	No	No	Yes	No	Yes	Bad
Catwoman	Female	Yes	No	No	Yes	No	Bad
Joker	Male	No	No	No	No	No	Bad

Classification: Good vs. Bad

➤ Predict a new hero (unlabeled data) as good or bad.



Test data

	<i>Gender</i>	<i>Mask</i>	<i>Cape</i>	<i>Tie</i>	<i>Ears</i>	<i>Smokes</i>	<i>Label</i>
Batgirl	Female	Yes	Yes	Yes	Yes	No	??
Riddler	Male	Yes	No	No	No	No	??

How to Learn the Decision Tree?



- Finding an optimal decision tree that correctly classifies a training set is an **NP-complete problem**.
 - ◆ Note: the optimal decision tree means the **smallest decision tree**.

- **Use a greedy heuristic!**
 - ◆ Split on the best attribute at each stage.



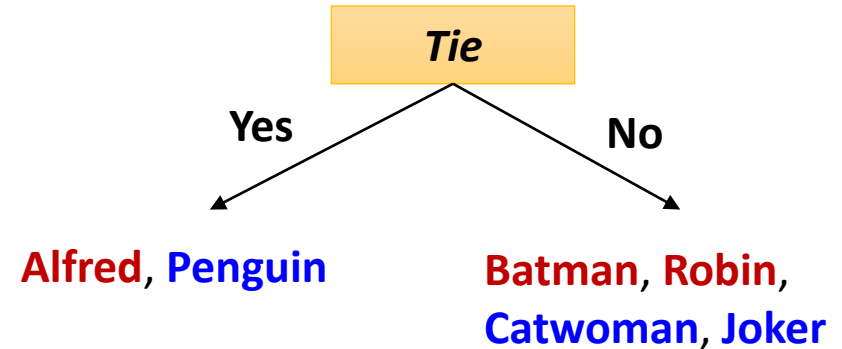
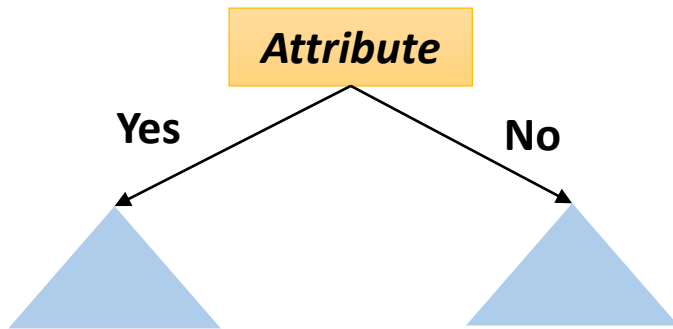
- **Which attribute is the best?**

- **When should we stop?**



Choosing a Good Split

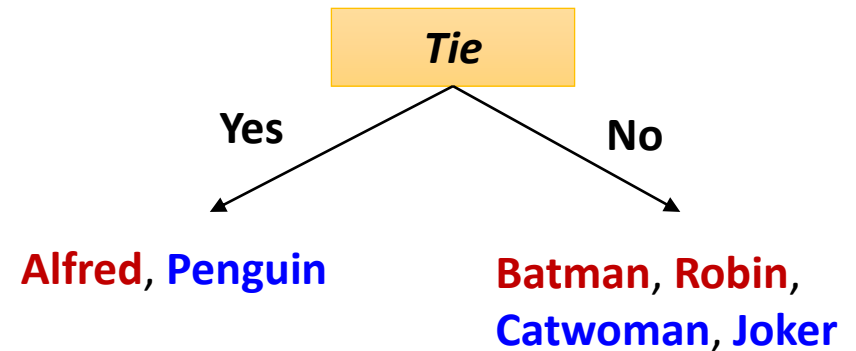
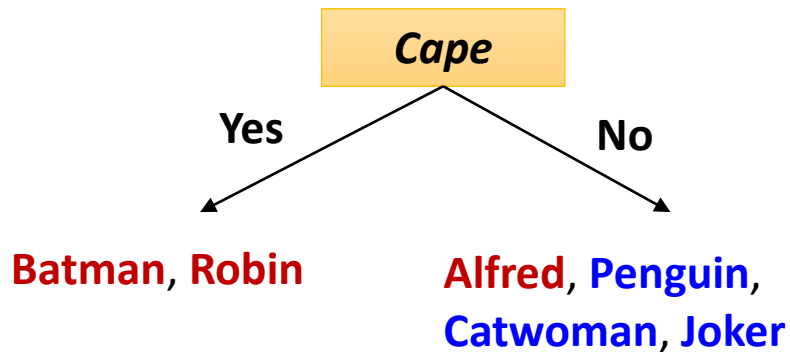
- Choosing a decision rule to split data into disjoint subsets



	<i>Gender</i>	<i>Mask</i>	<i>Cape</i>	<i>Tie</i>	<i>Ears</i>	<i>Smokes</i>	<i>Label</i>
Batman	Male	Yes	Yes	No	Yes	No	Good
Robin	Male	Yes	Yes	No	No	No	Good
Alfred	male	No	No	Yes	No	No	Good
Penguin	Male	No	No	Yes	No	Yes	Bad
Catwoman	Female	Yes	No	No	Yes	No	Bad
Joker	Male	No	No	No	No	No	Bad

Choosing a Good Split

➤ Which attribute is better?

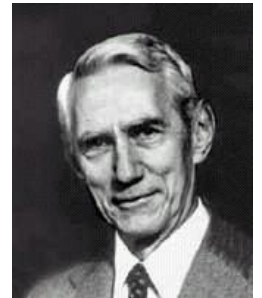


Choosing a Good Split



- How do we measure **uncertainty** in prediction for a leaf node?
 - ◆ All samples in the leaf node have the same class: **good**, low uncertainty
 - ◆ Each class has the same samples in the leaf node: **bad**, high uncertainty
- Idea: Define the **probability distribution** and use **information theory** to measure uncertainty.

A photograph of a chalkboard with the mathematical formula for entropy written in white chalk. The formula is $H = -\sum p(x) \log p(x)$. The chalkboard has a dark surface and a wooden frame.



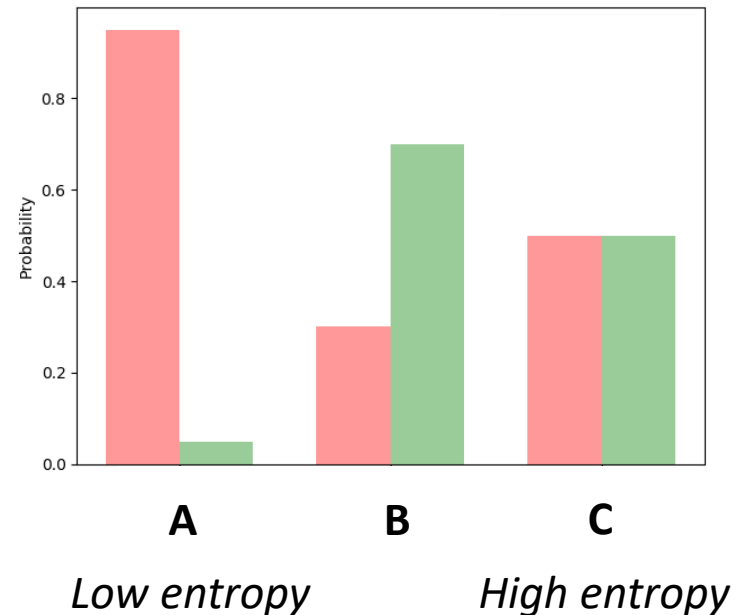
Information Entropy

- Given a discrete random variable X , information entropy is defined as follows.

$$H(X) = \mathbb{E}_{x \sim P(x)}[-\log_2 p(x)] = - \sum_{x \in X} p(x) \log_2 p(x)$$

➤ **Two extreme cases**

- ◆ Samples only have one class.
- ◆ Samples are divided into each class.



Information Entropy



➤ Measuring the degree of uncertainty

$$H(X) = -P_1 \log_2 P_1 - P_0 \log_2 P_0$$

P_1 = The probability that 1 appears in the set

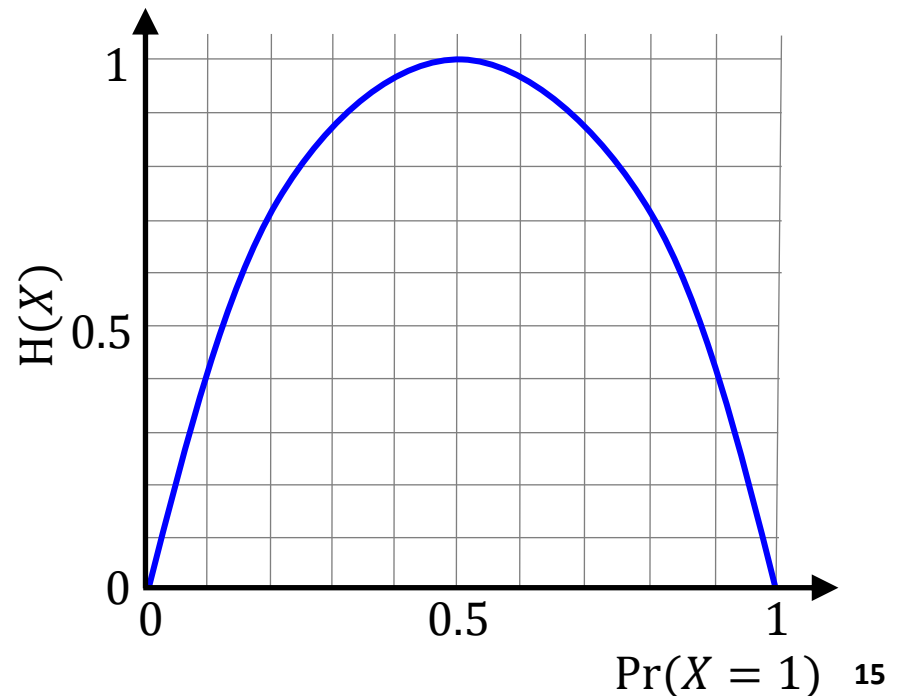
P_0 = The probability that 0 appears in the set

$$A = \{1, 1, 0, 0, 0, 0, 0, 0\}$$

$$H(X) = -2/8 \log_2 2/8 - 6/8 \log_2 6/8$$

$$B = \{1, 1, 1, 1, 0, 0, 0, 0\}$$

$$H(X) = -4/8 \log_2 4/8 - 4/8 \log_2 4/8$$



Information Entropy

- Assume that there are two classes ***P*** and ***N***.
- A dataset \mathcal{D} contains ***p*** elements of class ***P*** and ***n*** elements of class ***N***, respectively.
- For the dataset \mathcal{D} , the information entropy is computed by:

$$I(p, n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

Conditional Entropy

- Measuring **the average of the entropy** for each partition

$$H(Y | X) = \mathbb{E}_{x \sim P(x)} [H(Y | X = x)] = \sum_{x \in X} p(x) H(Y | X = x)$$

- Which partition is better?

{1,1,1,1, 0,0,0,0}



$A = \{\{1,1,0,0,0\}, \{1,1,0\}\}$

vs.

$B = \{\{1,1,1,1,0\}, \{0,0,0\}\}$

Conditional Entropy

$$A = \{\{1,1,0,0,0\}, \{1,1,0\}\}$$

$$-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$$

$$-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.92$$

$$\frac{5}{8} \times 0.97 + \frac{3}{8} \times 0.92 = 0.95$$

$$B = \{\{1,1,1,1,0\}, \{0,0,0\}\}$$

$$-\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} = 0.72$$

$$-\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3} = 0.00$$

$$\frac{5}{8} \times 0.72 + \frac{3}{8} \times 0.0 = 0.45$$

Information Gain

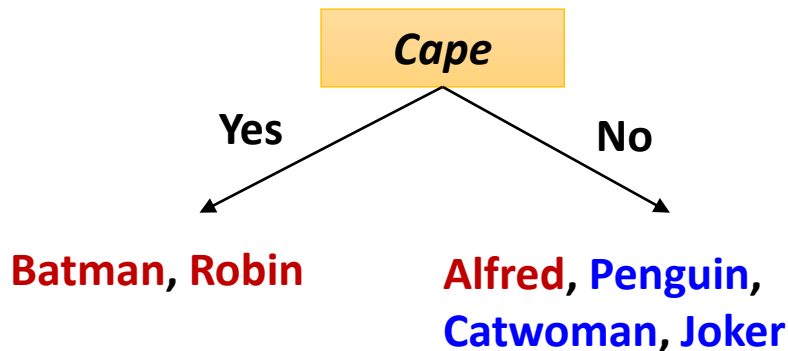
- The **information gain** in Y due to A , or the **mutual information** of Y and A is defined as follows.

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

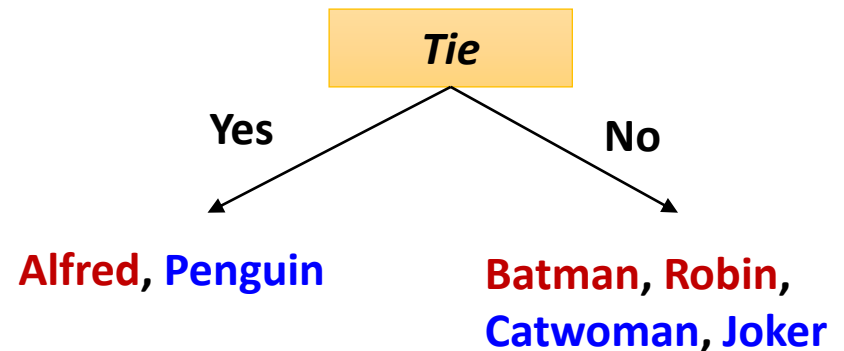
- If A is completely **informative** about Y : $IG(Y | A) = H(Y)$
 - ◆ The conditional entropy $H(Y | A)$ is 0.
- If A is completely **uninformative** about Y : $IG(Y | A) = 0$
 - ◆ The conditional entropy $H(Y | A)$ is equal to $H(Y)$.

Information Gain

➤ Suppose that an attribute A is chosen.



$$H(Y \mid \text{Cape}) = \frac{2}{6} \times I(1,0) + \frac{4}{6} \times I(3,1)$$

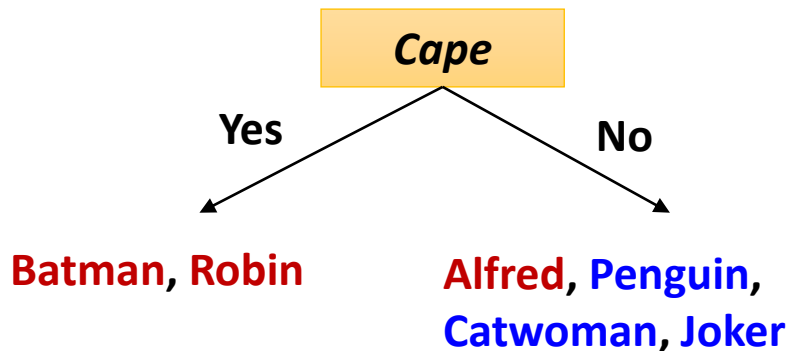


$$H(Y \mid \text{Tie}) = \frac{2}{6} \times I(1,1) + \frac{4}{6} \times I(2,2)$$

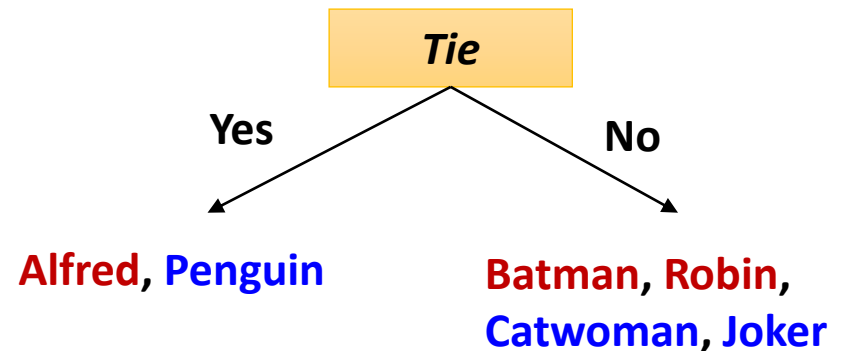
Information Gain



➤ Entropy before branching – entropy after branching



$$\begin{aligned} IG(Y | \text{Cape}) &= H(Y) - H(Y | \text{Cape}) \\ &= I(3, 3) - \left(\frac{2}{6} \times I(2, 0) + \frac{4}{6} \times I(1, 3) \right) \end{aligned}$$



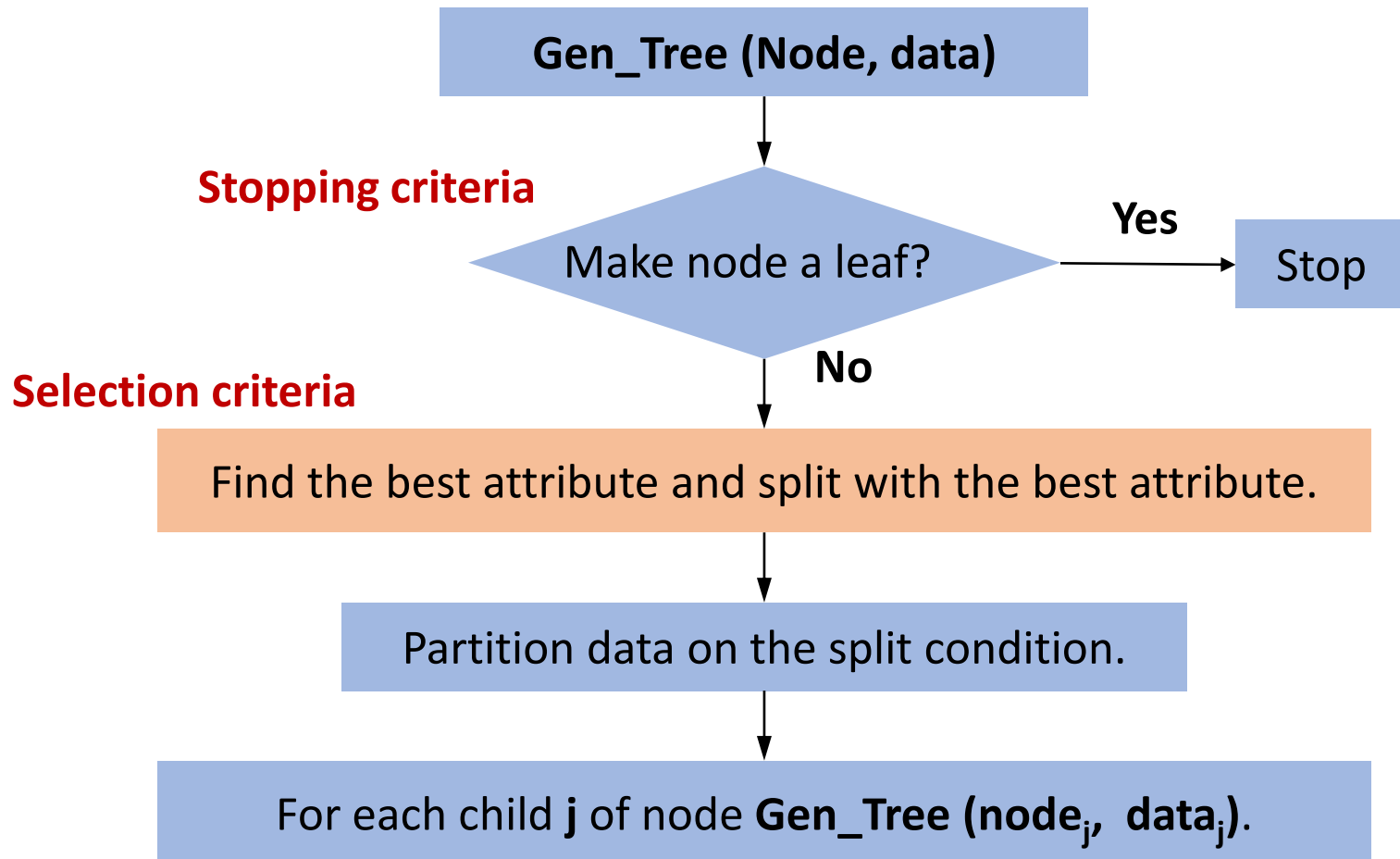
$$\begin{aligned} IG(Y | \text{Cape}) &= H(Y) - H(Y | \text{Tie}) \\ &= I(3, 3) - \left(\frac{2}{6} \times I(1, 1) + \frac{4}{6} \times I(2, 2) \right) \end{aligned}$$



Iterative Dichotomiser 3 (ID3)

Decision Trees (ID3)

- Top-down and **divide-and-conquer** approach





Stopping Condition

- All samples for a given node belong to **the same class**.
- There are **no samples left**.
- There are **no attributes** for additional partitioning.
 - ◆ The **majority voting** is used to determine the leaf node.

Training Procedure



➤ Selecting the attribute with the highest information gain

$$H(Y) = I(3, 3) = -\frac{3}{6}\log_2\frac{3}{6} - \frac{3}{6}\log_2\frac{3}{6} = 1.0$$

$$H(Y | \text{Gender}) = \frac{5}{6}I(3, 2) + \frac{1}{6}I(0, 1) = \frac{5}{6}\left(-\frac{3}{5}\log_2\frac{3}{5} - \frac{2}{5}\log_2\frac{2}{5}\right) = 0.809$$

$$\left. \begin{array}{l} H(Y) = 1.0 \\ H(Y | \text{Gender}) = 0.809 \end{array} \right\} \begin{array}{l} IG(Y | \text{Gender}) \\ = 1.0 - 0.809 \\ = 0.191 \end{array}$$

Training data

	<i>Gender</i>	<i>Mask</i>	<i>Cape</i>	<i>Tie</i>	<i>Ears</i>	<i>Smokes</i>	<i>Label</i>
Batman	Male	Yes	Yes	No	Yes	No	Good
Robin	Male	Yes	Yes	No	No	No	Good
Alfred	Male	No	No	Yes	No	No	Good
Penguin	Male	No	No	Yes	No	Yes	Bad
Catwoman	Female	Yes	No	No	Yes	No	Bad
Joker	Male	No	No	No	No	No	Bad

Training Procedure



➤ Selecting the attribute with the highest information gain

$$H(Y) = I(3, 3) = -\frac{3}{6}\log_2\frac{3}{6} - \frac{3}{6}\log_2\frac{3}{6} = 1.0$$

$$H(Y | Gender) = \frac{5}{6}I(3, 2) + \frac{1}{6}I(0, 1) = \frac{5}{6}\left(-\frac{3}{5}\log_2\frac{3}{5} - \frac{2}{5}\log_2\frac{2}{5}\right) = 0.809$$

$$H(Y | Mask) = \frac{3}{6}I(2, 1) + \frac{3}{6}I(1, 2) = -\frac{2}{3}\log_2\frac{2}{3} - \frac{1}{3}\log_2\frac{1}{3} = 0.918$$

$$H(Y | Cape) = \frac{2}{6}I(2, 0) + \frac{4}{6}I(1, 3) = \frac{4}{6}\left(-\frac{1}{4}\log_2\frac{1}{4} - \frac{3}{4}\log_2\frac{3}{4}\right) = 0.540$$

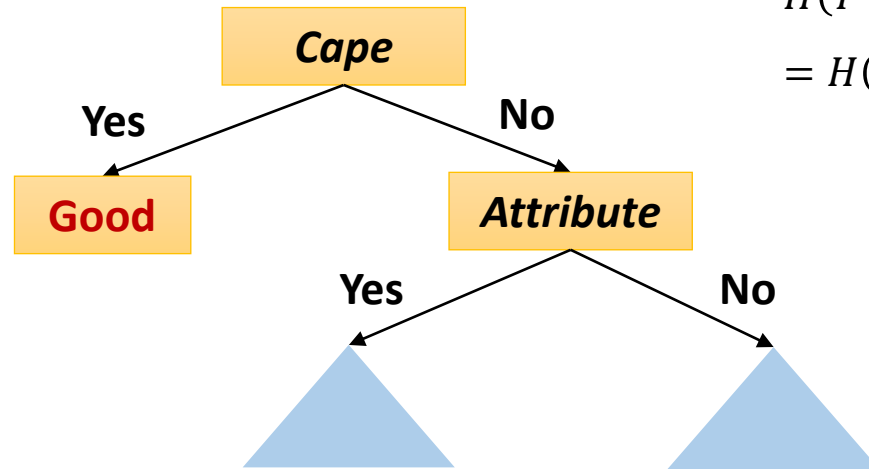
$$H(Y | Tie) = \frac{2}{6}I(1, 1) + \frac{4}{6}I(2, 2) = 1.0$$

$$H(Y | Ears) = \frac{2}{6}I(1, 1) + \frac{4}{6}I(2, 2) = 1.0$$

$$H(Y | Smokes) = \frac{1}{6}I(0, 1) + \frac{5}{6}I(3, 2) = 0.809$$

Training Procedure

- For the subset (*Cape* = *no*), select the attribute with the highest information gain.



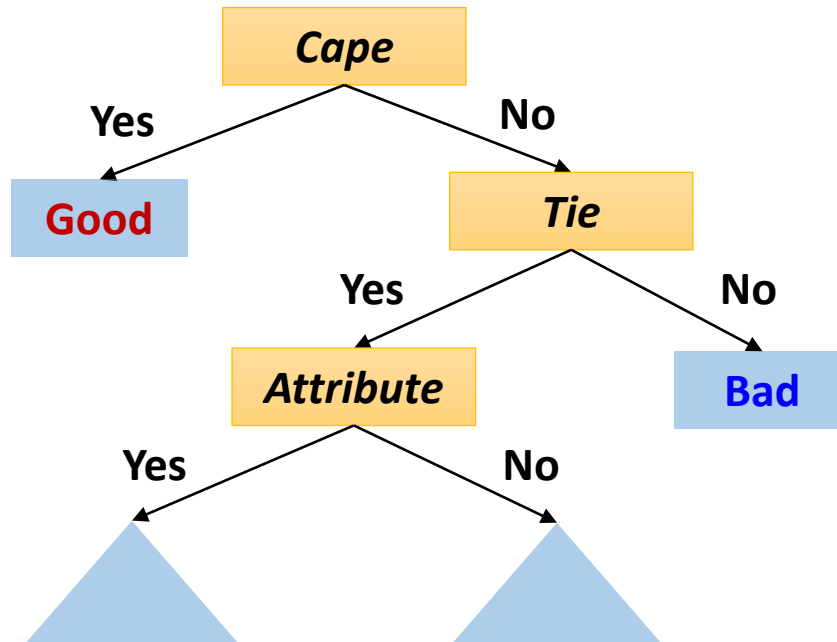
$$\begin{aligned}
 H(Y \mid \text{Gender}) &= H(Y \mid \text{Mask}) = H(Y \mid \text{Ears}) \\
 &= H(Y \mid \text{Smokes}) = \frac{3}{4}I(1, 2) + \frac{1}{4}I(0, 1) = 0.689
 \end{aligned}$$

$$H(Y \mid \text{Tie}) = \frac{2}{4}I(1, 1) + \frac{2}{4}I(0, 2) = 0.5$$

	<i>Gender</i>	<i>Mask</i>	<i>Tie</i>	<i>Ears</i>	<i>Smokes</i>	<i>Label</i>
Alfred	Male	No	Yes	No	No	Good
Penguin	Male	No	Yes	No	Yes	Bad
Catwoman	Female	Yes	No	Yes	No	Bad
Joker	Male	No	No	No	No	Bad

Training Procedure

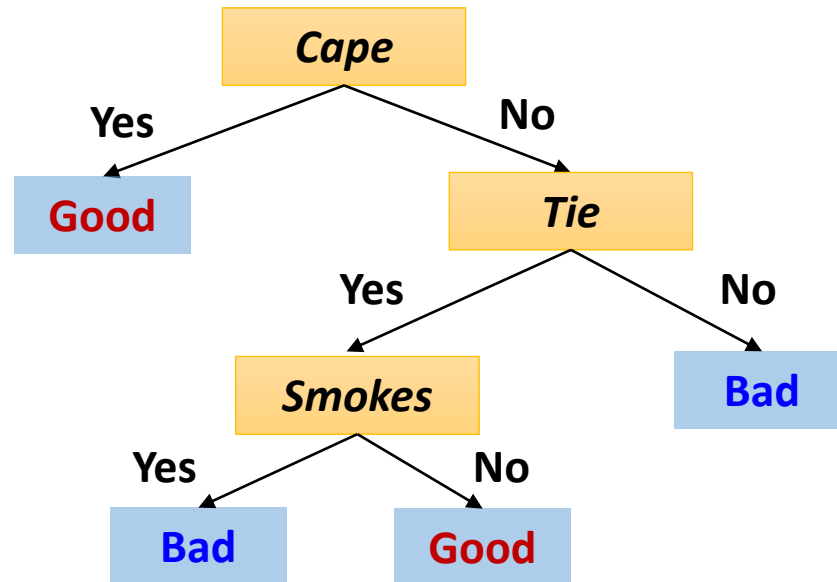
- For the subset ($Tie = yes$), select the attribute with the highest information gain.



	<i>Gender</i>	<i>Mask</i>	<i>Ears</i>	<i>Smokes</i>	<i>Label</i>
Alfred	Male	No	No	No	Good
Penguin	Male	No	No	Yes	Bad

Making Prediction

➤ Predicting the label for each sample on test data



Testing data

	<i>Gender</i>	<i>Mask</i>	<i>Cape</i>	<i>Tie</i>	<i>Ears</i>	<i>Smokes</i>	<i>Label</i>
Batgirl	Female	Yes	Yes	Yes	Yes	No	??
Riddler	Male	Yes	No	No	No	No	??

Example: Buying Laptops



Training data

<i>Age</i>	<i>Income</i>	<i>Student</i>	<i>Credit</i>	<i>Buy</i>
<= 30	High	N	Fair	No
<= 30	High	N	Excellent	No
31 ... 40	High	N	Fair	Yes
> 40	Medium	N	Fair	Yes
> 40	Low	Y	Fair	Yes
> 40	Low	Y	Excellent	No
31 ... 40	Low	Y	Excellent	Yes
<= 30	Medium	N	Fair	No
<= 30	Low	Y	Fair	Yes
> 40	Medium	Y	Fair	Yes
<= 30	Medium	Y	Excellent	Yes
31 ... 40	Medium	N	Excellent	Yes
31 ... 40	High	Y	Fair	Yes
> 40	Medium	N	Excellent	No

Example: Buying Laptops

➤ Selecting the attribute with the highest information gain

<i>Age</i>	<i>p</i>	<i>n</i>	<i>I(p, n)</i>
<= 30	2	3	0.971
30 ... 40	4	0	0
> 40	3	2	0.971

<i>Age</i>	<i>Buy</i>
<= 30	No
<= 30	No
31 ... 40	Yes
> 40	Yes
> 40	Yes
> 40	No
31 ... 40	Yes
<= 30	No
<= 30	Yes
> 40	Yes
<= 30	Yes
31 ... 40	Yes
31 ... 40	Yes
> 40	No

$$\begin{aligned}
 & H(Y) - H(Y | Age) \\
 &= H(Y) - \left(\frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(2,3) \right) = \\
 &= 0.940 - 0.69 = 0.25
 \end{aligned}$$

Example: Buying Laptops

➤ Selecting the attribute with the highest information gain

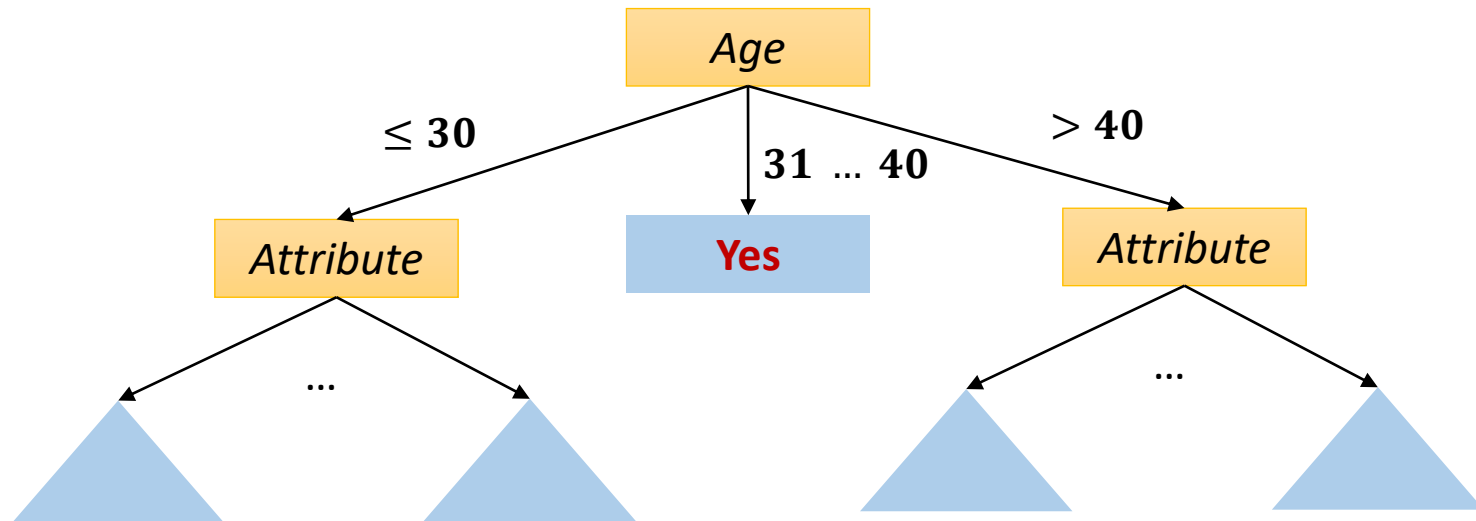
<i>Student</i>	<i>p</i>	<i>n</i>	<i>I(p, n)</i>
Yes	6	1	0.971
No	4	3	0

$$\begin{aligned}
 & H(Y) - H(Y \mid \text{Student}) \\
 &= H(Y) - \left(\frac{7}{14} I(6,1) + \frac{7}{14} I(4,3) \right) = 0.151
 \end{aligned}$$

<i>Student</i>	<i>Buy</i>
N	No
N	No
N	Yes
N	Yes
Y	Yes
Y	No
Y	Yes
N	No
Y	Yes
Y	Yes
Y	Yes
N	Yes
Y	Yes
N	No

Example: Buying Laptops

➤ A decision tree after the first partitioning



Income	Student	Credit	Buy
High	N	Fair	No
High	N	Excellent	No
Medium	N	Fair	No
Low	Y	Fair	Yes
Medium	Y	Excellent	Yes

Income	Student	Credit	Buy
Medium	N	Fair	Yes
Low	Y	Fair	Yes
Low	Y	Excellent	No
Medium	Y	Fair	Yes
Medium	N	Excellent	No

Example: Buying Laptops

➤ Using age, it is partitioned into three subsets.

- ◆ For each subset, select the attribute with the highest information gain in a recursive manner.

<i>Income</i>	<i>Student</i>	<i>Credit</i>	<i>Buy</i>
High	N	Fair	No
High	N	Excellent	No
Medium	N	Fair	No
Low	Y	Fair	Yes
Medium	Y	Excellent	Yes

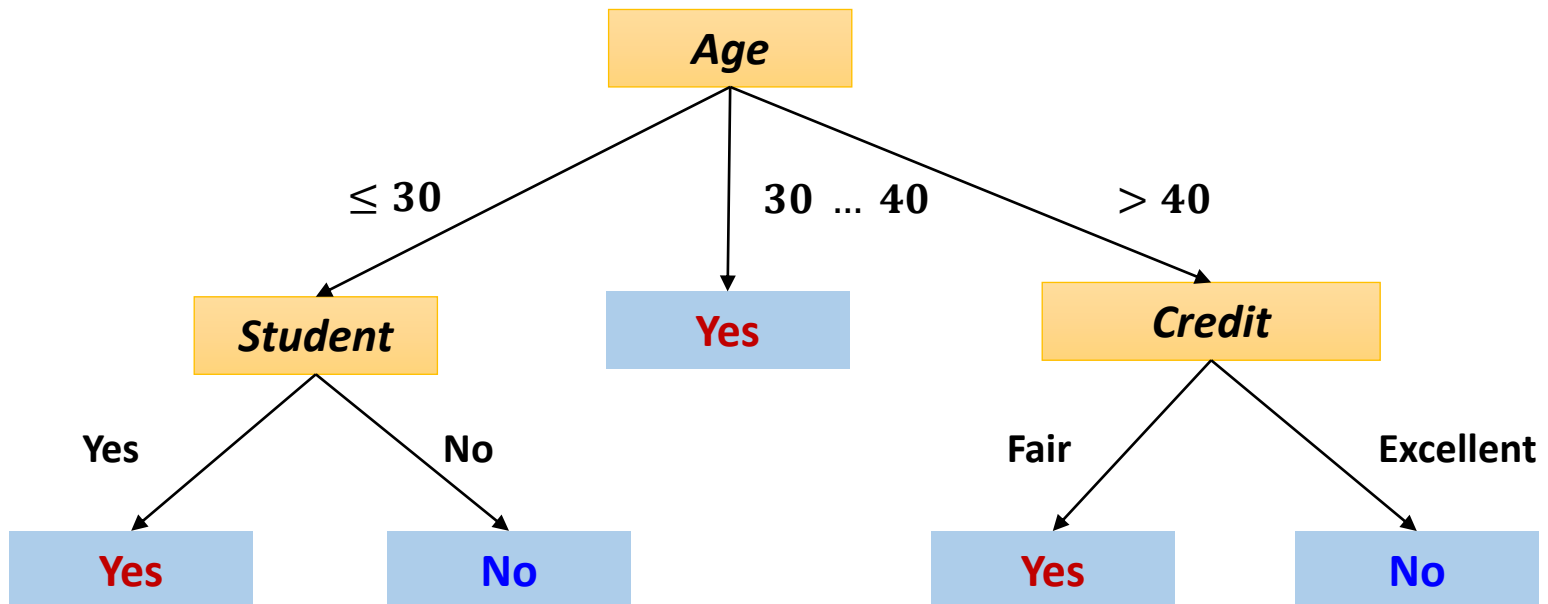
➤ Let each subset \mathcal{D}_1 , \mathcal{D}_2 , and \mathcal{D}_3 .

- ◆ \mathcal{D}_2 has already been decided as a single label.

<i>Income</i>	<i>Student</i>	<i>Credit</i>	<i>Buy</i>
High	N	Fair	Yes
Low	Y	Excellent	Yes
Medium	N	Excellent	Yes
High	Y	Fair	Yes

<i>Income</i>	<i>Student</i>	<i>Credit</i>	<i>Buy</i>
Medium	N	Fair	Yes
Low	Y	Fair	Yes
Low	Y	Excellent	No
Medium	Y	Fair	Yes
Medium	N	Excellent	No

Example: Buying Laptops





C4.5: Improving ID3



C4.5: Improving ID3

- C4.5 has several improvements to ID3.
- Handling both **continuous** and **discrete** attributes
- Handling **attributes with differing costs**
- Pruning trees after creation

Handling Continuous Attributes



Training data

<i>Outlook</i>	<i>Temperature</i>	<i>Humidity</i>	<i>Wind</i>	<i>Play</i>
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	78	False	Yes
Rainy	70	96	False	Yes
Rainy	68	80	False	Yes
Rainy	65	70	True	No
Overcast	64	65	True	Yes
Sunny	72	95	False	No
Sunny	69	70	False	Yes
Rainy	75	80	False	Yes
Sunny	75	70	True	Yse
Overcast	72	90	True	Yse
Overcast	81	75	False	Yse
Rainy	71	80	True	No

Handling Continuous Attributes



➤ Finding the best split point by **enumerating all possible cases**

Temp	Humidity	Wind	Play
85	85	False	No
80	90	True	No
83	78	False	Yes
70	96	False	Yes
68	80	False	Yes
65	70	True	No
64	65	True	Yes
72	95	False	No
69	70	False	Yes
75	80	False	Yes
75	70	True	Yes
72	90	True	Yes
81	75	False	Yes
71	80	True	No

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

Handling Continuous Attributes



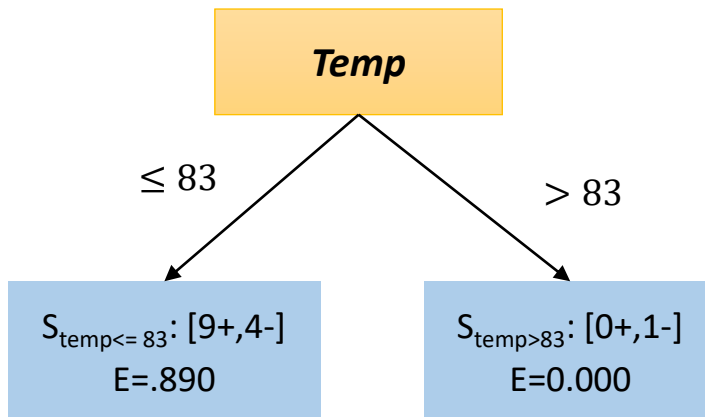
➤ Finding the best split point by **enumerating all possible cases**

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

Handling Continuous Attributes

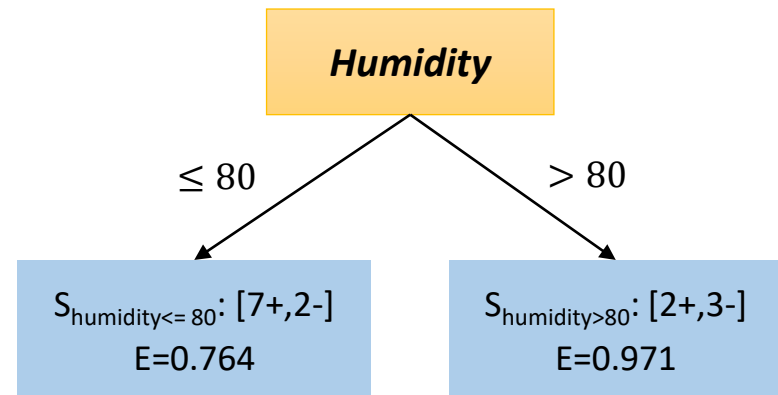


$S: [9+, 5-], E = 0.940$



$$\begin{aligned} \text{Gain}(temp) &= .940 - (13/14) \times .890 - (1/14) \\ &\times 1.0 = .113 \end{aligned}$$

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



$$\begin{aligned} \text{Gain}(humidity) &= .940 - (9/14) \times 0.764 - (5/14) \\ &\times .971 = .102 \end{aligned}$$

Handling Continuous Attributes



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

outlook

Sunny

Overcast

Rain

$S_{\text{outlook=Sunny}}: [2+, 3-]$

$E=0.971$

$S_{\text{outlook=Overcast}}: [4+, 0-]$

$E=0.00$

$S_{\text{outlook=Rain}}: [3+, 2-]$

$E=.971$

$\text{Gain}(\text{outlook})$

$$\begin{aligned} &= .940 - (5/14) \times .971 - (4/14) \times 1.0 \\ &\quad - (5/14) \times .971 = .246 \end{aligned}$$

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

wind

Weak

Strong

$S_{\text{wind=Weak}}: [6+, 2-]$

$E=0.811$

$S_{\text{wind=Strong}}: [3+, 3-]$

$E=1.00$

$\text{Gain}(\text{wind})$

$$\begin{aligned} &= .940 - (8/14) \times .811 - (6/14) \\ &\quad \times 1.0 = .048 \end{aligned}$$

Problem of Information Gain

➤ Selecting the attribute with the highest information gain

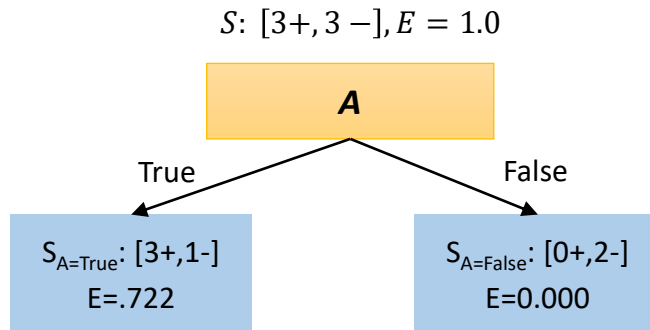
<i>A</i>	<i>B</i>	<i>Play</i>
True	B1	Y
True	B2	Y
True	B3	Y
True	B4	N
False	B5	N
False	B6	N

➤ Which attribute is better?

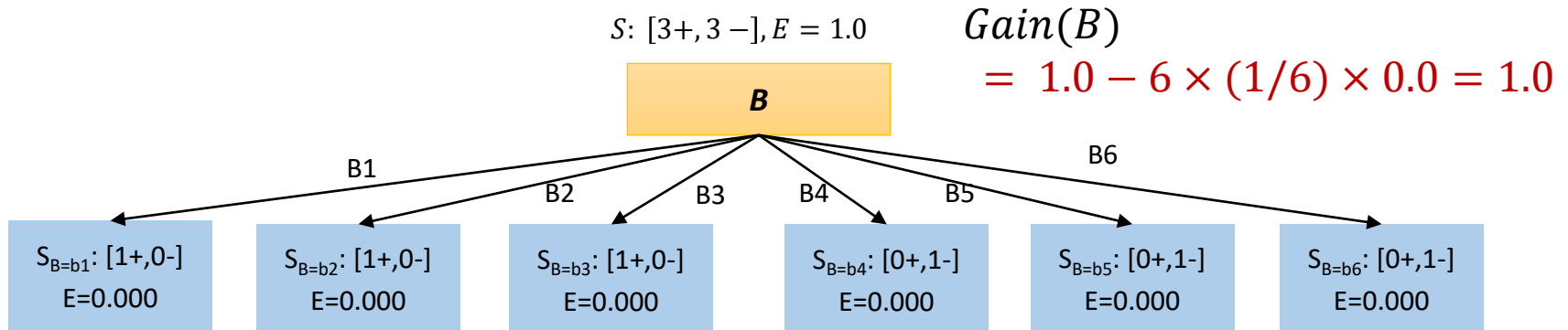
Problem of Information Gain



- Entropy does not consider a **branching factor**.



$$\begin{aligned}
 \text{Gain}(A) &= 1.0 - (4/6) \times .722 - (2/6) \times 1.0 \\
 &= .484
 \end{aligned}$$



$$\begin{aligned}
 \text{Gain}(B) &= 1.0 - 6 \times (1/6) \times 0.0 = 1.0
 \end{aligned}$$

GainInfo: Alternative Measure

➤ How to considering a **branching factor**?

$$SplitInfo(Y) = - \sum_{i=1}^k \frac{|\mathcal{D}_i|}{|\mathcal{D}|} \log_2 \frac{|\mathcal{D}_i|}{|\mathcal{D}|}$$

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k\}$$

A	B	Play
True	B1	Y
True	B2	Y
True	B3	Y
True	B4	N
False	B5	N
False	B6	N

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

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

GainInfo: Alternative Measure

➤ How to considering a **branching factor**?

$$GainInfo(Y) = \frac{Gain(Y)}{SplitInfo(Y)}$$

A	B	Play
True	B1	Y
True	B2	Y
True	B3	Y
True	B4	N
False	B5	N
False	B6	N

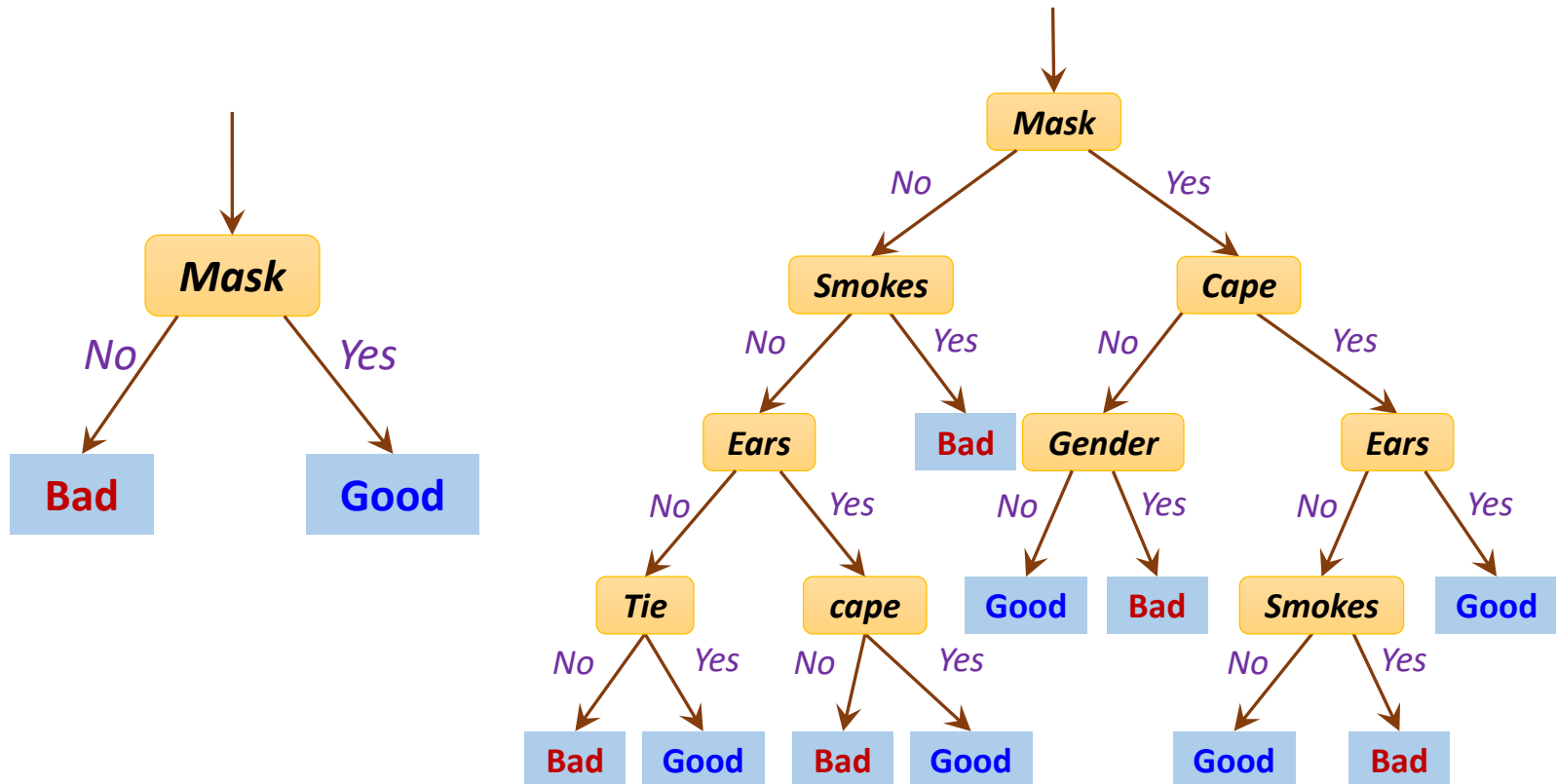
$$GainInfo(A) = \frac{0.484}{0.918} = 0.527$$

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

➤ Therefore, A is better than B.

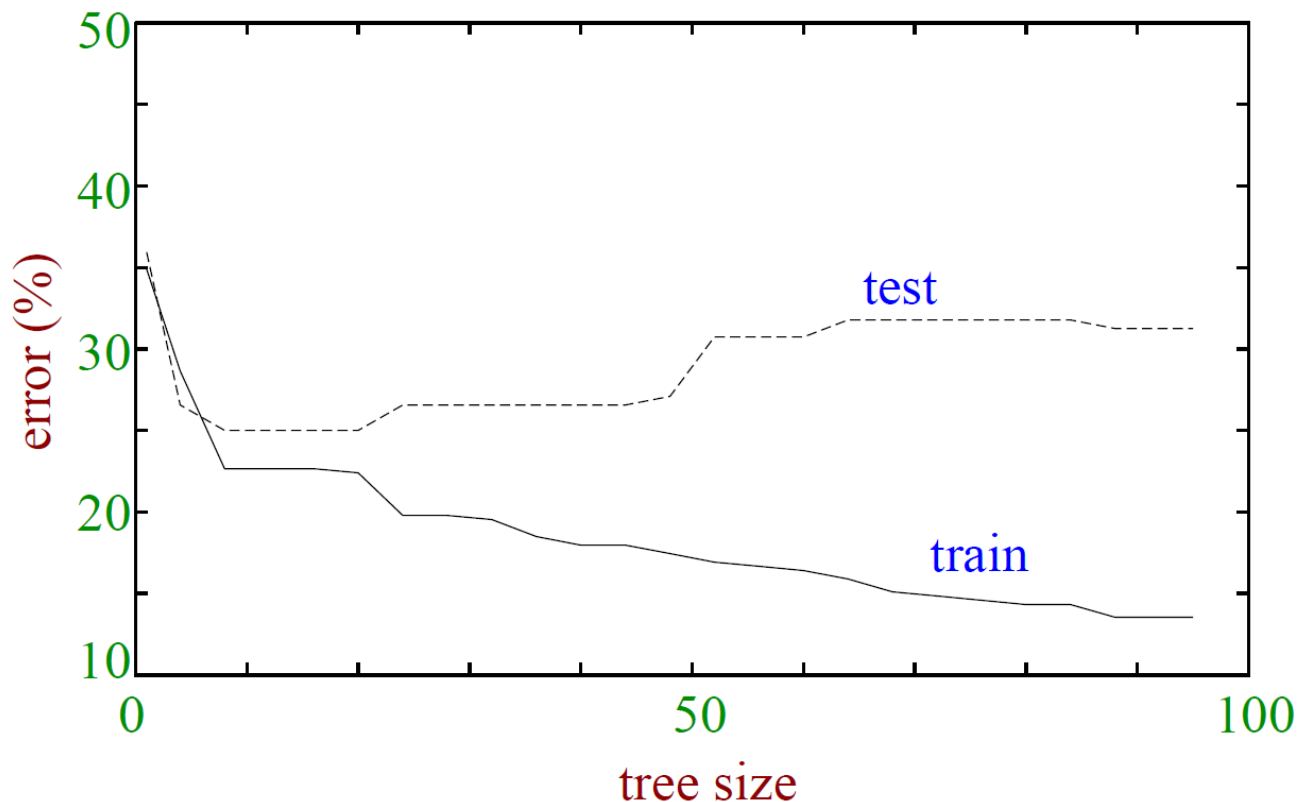
Underfitting vs. Overfitting

- Underfitting: overly simple, does not even fit available data
- Overfitting: too complicated, perfectly classifies training data



Tree Size vs. Accuracy

- Trees must be big enough to fit training data.
- However, **too big trees may overfit.**
 - ◆ It is difficult to decide the **best tree size** from training errors.



Two Solutions to Avoid Overfitting

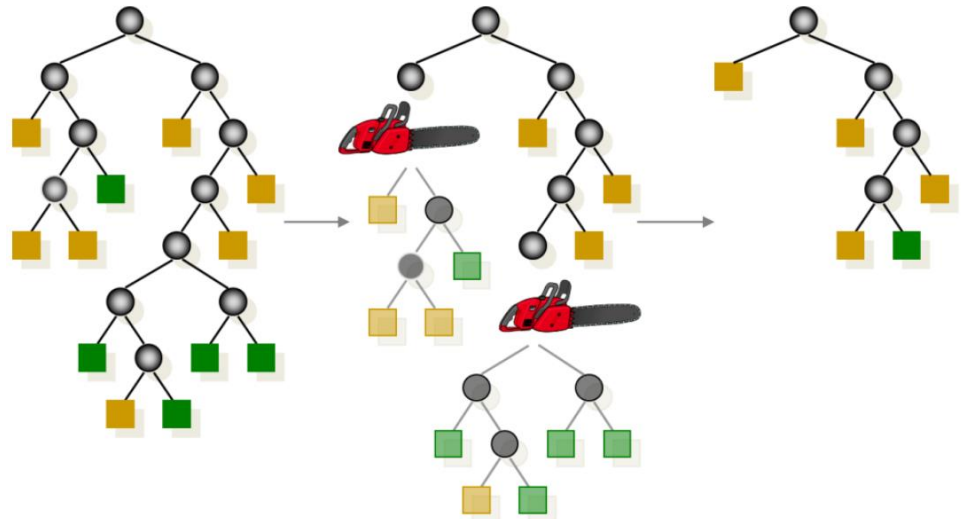


➤ Prepruning

- ◆ Do not split a node if this would result in **the goodness measure falling below a threshold.**
- ◆ It is difficult to choose an **appropriate threshold.**

➤ Postpruning

- ◆ Remove **branches** from a **fully grown tree.**



Q&A

