

# Data Mining, Data Warehouse

## Decision Tree Learning

### Lesson #4

# Attribute Selection Measure in ID3

## Information Gain

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} \cdot Entropy(S_v)$$

Select the attribute  $A$  that best classifies  
the *examples*

# ...An Example

Weekend (Example)	Weather	Parents	Money	Decision (Class)
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay in
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

# Using the ID3 algorithm to build a decision tree

- Figure out which attribute will be put into the node at the top of our tree:

*Weather*, *parents* or *money*.

- To do this, we first need to calculate:

Entropy(S) =

$$-p_{\text{cinema}} \log_2(p_{\text{cinema}}) - p_{\text{tennis}} \log_2(p_{\text{tennis}}) - p_{\text{shopping}} \log_2(p_{\text{shopping}}) - p_{\text{stay\_in}} \log_2(p_{\text{stay\_in}})$$

$$\begin{aligned} &= -(6/10) * \log_2(6/10) - (2/10) * \log_2(2/10) - (1/10) * \log_2(1/10) - (1/10) * \log_2(1/10) \\ &= 1.571 \end{aligned}$$

Table



# Using the ID3 algorithm... (cont.)

Now let's calculate the gain for the *weather* attribute:

$$\text{Gain}(S, \text{weather}) = 1.571 - (|S_{\text{sunny}}|/10) * \text{Entropy}(S_{\text{sunny}}) - (|S_{\text{windy}}|/10) * \text{Entropy}(S_{\text{windy}}) - (|S_{\text{rainy}}|/10) * \text{Entropy}(S_{\text{rainy}})$$

**Where,**

$$\begin{aligned} \text{Entropy}(S_{\text{sunny}}) &= -p_{\text{cinema}|\text{sunny}} \log_2(p_{\text{cinema}|\text{sunny}}) - p_{\text{tennis}|\text{sunny}} \log_2(p_{\text{tennis}|\text{sunny}}) \\ &= -(1/3) * \log_2(1/3) - (2/3) * \log_2(2/3) = 0.918 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(S_{\text{windy}}) &= -p_{\text{cinema}|\text{windy}} \log_2(p_{\text{cinema}|\text{windy}}) \\ &\quad - p_{\text{shopping}|\text{windy}} \log_2(p_{\text{shopping}|\text{windy}}) \\ &= -(1/4) * \log_2(1/4) - (3/4) * \log_2(3/4) = 0.811 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(S_{\text{rainy}}) &= -p_{\text{cinema}|\text{rainy}} \log_2(p_{\text{cinema}|\text{rainy}}) - p_{\text{stay\_in}|\text{rainy}} \log_2(p_{\text{stay\_in}|\text{rainy}}) \\ &= -(2/3) * \log_2(2/3) - (1/3) * \log_2(1/3) = 0.918 \end{aligned}$$

$$\text{Gain}(S, \text{weather}) = 1.571 - (3/10) * \text{Entropy}(S_{\text{sunny}}) - (4/10) * \text{Entropy}(S_{\text{windy}}) - (3/10) * \text{Entropy}(S_{\text{rainy}}) = \mathbf{0.69}$$

Table

# Using the ID3 algorithm to build a decision tree (cont.)

Now let's calculate the gain for the *parents* attribute:

$$\text{Gain}(S, \text{parents}) = 1.571 - (|S_{\text{yes}}|/10) * \text{Entropy}(S_{\text{yes}}) - (|S_{\text{no}}|/10) * \text{Entropy}(S_{\text{no}})$$

**Where,**

$$\begin{aligned} \text{Entropy}(S_{\text{yes}}) &= -p_{\text{cinema}|\text{yes}} \log_2(p_{\text{cinema}|\text{yes}}) \\ &= -(5/5) * \log_2(5/5) = 0 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(S_{\text{no}}) &= -p_{\text{tennis}|\text{no}} \log_2(p_{\text{tennis}|\text{no}}) - p_{\text{stay\_in}|\text{no}} \log_2(p_{\text{stay\_in}|\text{no}}) \\ &\quad - p_{\text{cinema}|\text{no}} \log_2(p_{\text{cinema}|\text{no}}) - p_{\text{shopping}|\text{no}} \log_2(p_{\text{shopping}|\text{no}}) \\ &= -(2/5) * \log_2(2/5) - (1/5) * \log_2(1/5) - (1/5) * \log_2(1/5) - (1/5) * \log_2(1/5) = 1.922 \end{aligned}$$

$$\text{Gain}(S, \text{parents}) = 1.571 - (5/10) * \text{Entropy}(S_{\text{yes}}) - (5/10) * \text{Entropy}(S_{\text{no}}) = \mathbf{0.61}$$

# Using the ID3 algorithm... (cont.)

Now let's calculate the gain for the *money* attribute:

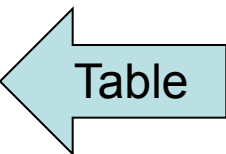
$$\text{Gain}(S, \text{money}) = 1.571 - (|S_{\text{rich}}|/10) * \text{Entropy}(S_{\text{rich}}) - (|S_{\text{poor}}|/10) * \text{Entropy}(S_{\text{poor}})$$

**Where,**

$$\begin{aligned} \text{Entropy}(S_{\text{rich}}) &= -p_{\text{cinema}|\text{rich}} \log_2(p_{\text{cinema}|\text{rich}}) - p_{\text{tennis}|\text{rich}} \log_2(p_{\text{tennis}|\text{rich}}) \\ &\quad - p_{\text{stay\_in}|\text{rich}} \log_2(p_{\text{stay\_in}|\text{rich}}) - p_{\text{shopping}} \log_2(p_{\text{shopping}}) \\ &= -(3/7) * \log_2(3/7) - (2/7) * \log_2(2/7) - (1/7) * \log_2(1/7) - (1/7) * \log_2(1/7) = 1.842 \end{aligned}$$

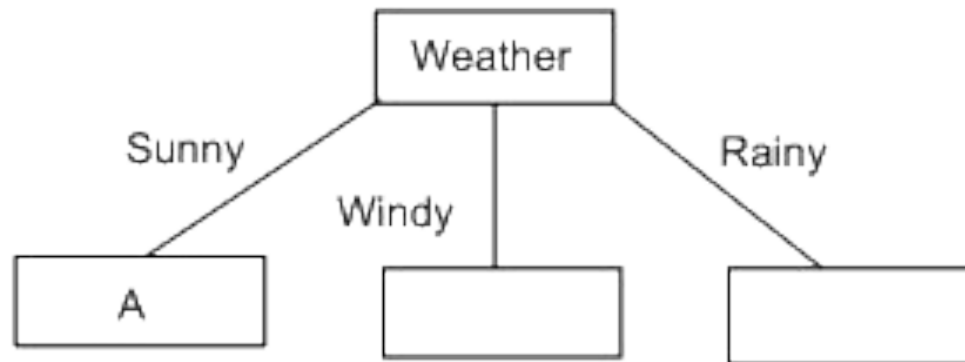
$$\begin{aligned} \text{Entropy}(S_{\text{poor}}) &= -p_{\text{cinema}|\text{no}} \log_2(p_{\text{cinema}|\text{no}}) \\ &= -(3/3) * \log_2(3/3) = 0 \end{aligned}$$

$$\text{Gain}(S, \text{money}) = 1.571 - (7/10) * \text{Entropy}(S_{\text{rich}}) - (3/10) * \text{Entropy}(S_{\text{poor}}) = \mathbf{0.28}$$



# Using the ID3 algorithm... (cont.)

- The first node in the decision tree will be the *weather* attribute.
- Now we look at the first branch.  $S_{\text{sunny}} = \{W1, W2, W10\}$ . This is not empty, so we do not put a default categorization leaf node here.





# Using the ID3 algorithm... (cont.)

- Now we have to fill in the choice of attribute A, which we know cannot be weather, because we've already removed that from the list of attributes to use.
- So, we need to calculate the values for  $\text{Gain}(S_{\text{sunny}}, \text{parents})$  and  $\text{Gain}(S_{\text{sunny}}, \text{money})$ .

Weekend (Example)	Weather	Parents	Money	Decision (Category)
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W10	Sunny	No	Rich	Tennis

# Using the ID3 algorithm... (cont.)

## ■ Hence we can calculate:

- $\text{Gain}(S_{\text{sunny}}, \text{parents}) = 0.918 - (|S_{\text{yes}}|/|S|) * \text{Entropy}(S_{\text{yes}}) - (|S_{\text{no}}|/|S|) * \text{Entropy}(S_{\text{no}})$   
 $= 0.918 - (1/3) * 0 - (2/3) * 0 = \mathbf{0.918}$

(Entropy( $S_{\text{sunny}}$ ))

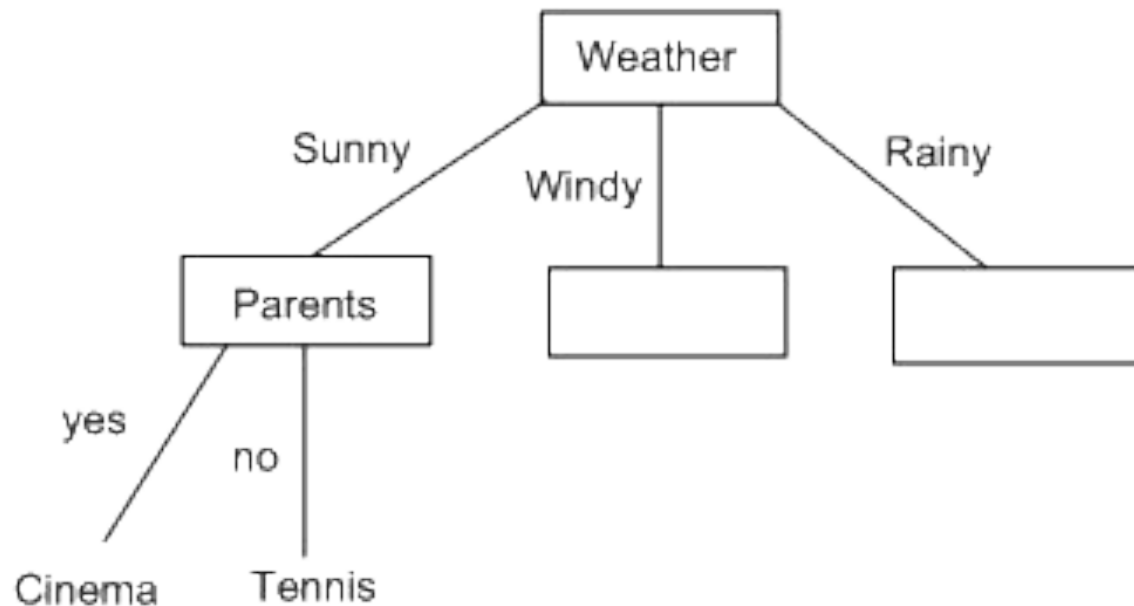
- $\text{Gain}(S_{\text{sunny}}, \text{money}) = 0.918 - (|S_{\text{rich}}|/|S|) * \text{Entropy}(S_{\text{rich}}) - (|S_{\text{poor}}|/|S|) * \text{Entropy}(S_{\text{poor}})$   
 $= 0.918 - (3/3) * 0.918 - 0 = 0.918 - 0.918 = 0$

$\log(1/3) * 1/3 - (2/3 * \log(2/3))$

- Note:  $\text{Entropy}(S_{\text{yes}})$  and  $\text{Entropy}(S_{\text{no}})$  were both zero, because  $S_{\text{yes}}$  contains examples which are all in the same category (cinema), and  $S_{\text{no}}$  contains examples which are all in the same category (tennis).
- This should make it more obvious why we use information gain to choose attributes to put in nodes.

# Using the ID3 algorithm... (cont.)

- Given our calculations, attribute A should be taken as *parents*.



# Using the ID3 algorithm... (cont.)

- Now we need to calculate the values for  $\text{Gain}(S_{\text{windy}}, \text{parents})$  and  $\text{Gain}(S_{\text{windy}}, \text{money})$ .

Weekend (Example)	Weather	Parents	Money	Decision (Category)
W3	Windy	Yes	Rich	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema

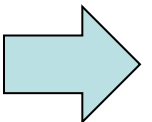
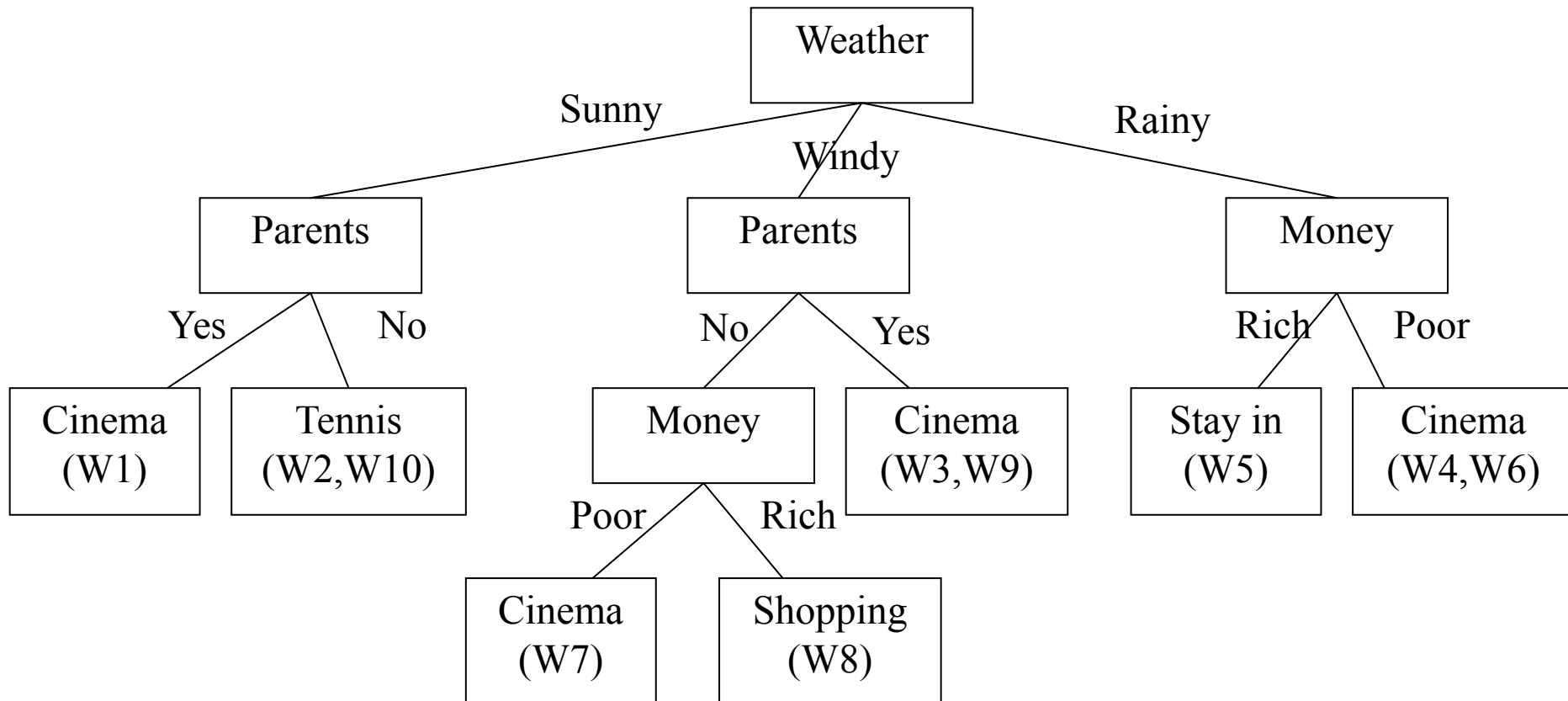
# Using the ID3 algorithm... (cont.)

- The calculation:
- $\text{Gain}(S_{\text{windy}}, \text{parents}) = 0.811 - (|S_{\text{yes}}|/|S|) * \text{Entropy}(S_{\text{yes}}) - (|S_{\text{no}}|/|S|) * \text{Entropy}(S_{\text{no}})$   
 $= 0.811 - (2/4) * 0 - (2/4) * 1 = \mathbf{0.311}$
- $\text{Gain}(S_{\text{windy}}, \text{money}) = 0.811 - (|S_{\text{rich}}|/|S|) * \text{Entropy}(S_{\text{rich}}) - (|S_{\text{poor}}|/|S|) * \text{Entropy}(S_{\text{poor}})$   
 $= 0.811 - (3/4) * 0.918 - (1/4) * 0 = 0.122$

Meaning that this node will be split by *parents* too.

After calculating the Gain for *rainy* too, the final tree is...

# Using the ID3 algorithm... (cont.)



# Accuracy Estimation

- **Training** Accuracy Rate
  - The percentage of *training set* samples that are correctly classified by the model.
- **Testing** Accuracy Rate
  - The percentage of *test set* samples that are correctly classified by the model.
  - Test set is independent of training set, otherwise over-fitting will occur.
- Majority Rule Accuracy
  - To select a class for a terminal node, select the class A having the most examples (in training set).
  - *Majority rule accuracy* =  $|S_A|/|S|$
  - The classification model should be more accurate than the *majority rule*.

# Accuracy Estimation

- In our example, training accuracy is  $10/10 = 100\%$
- Take the following testing set for example:

Weekend (Example)	Weather	Parents	Money	Decision (Class)
W11	Sunny	Yes	Poor	Cinema
W12	Sunny	No	Rich	Tennis
W13	Windy	No	Rich	Cinema
W14	Rainy	Yes	Poor	Cinema

Testing accuracy is:  $3/4 = 75\%$

Majority rule accuracy =  $6/10 = 60\%$



# Accuracy Estimation

- How can we tell if the **training accuracy** is significantly different than the **testing accuracy**?
- Normal Approximation to Binomial Distribution  
H0: The two accuracies are equal. We reject H0

if:  $\hat{p} > p + Z \cdot \sqrt{\frac{p(1-p)}{n}}$

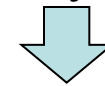
number of rows in  
the testing set

Testing accuracy

$$1 > p + 1.96 \cdot \sqrt{\frac{p(1-p)}{4}}$$

We can check in the same manner,  
majority rule vs. testing accuracy

If training accuracy > testing accuracy



Over-fitting

$p=0.51$  □ **Testing accuracy** should be 51% at the most, for the two accuracies to be significantly different (to reject H0).