

Learning Classification Trees

Sargur Srihari

srihari@cedar.buffalo.edu

Topics in CART

- CART as an adaptive basis function model
- Classification and Regression Tree Basics
- Growing a Tree

A Classification Tree

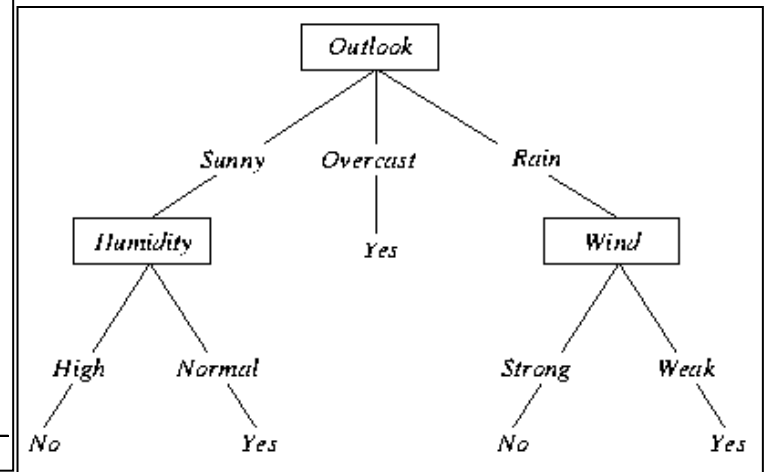
Data Set

$\mathbf{x} = [x_1, x_2, x_3, x_4]$

y

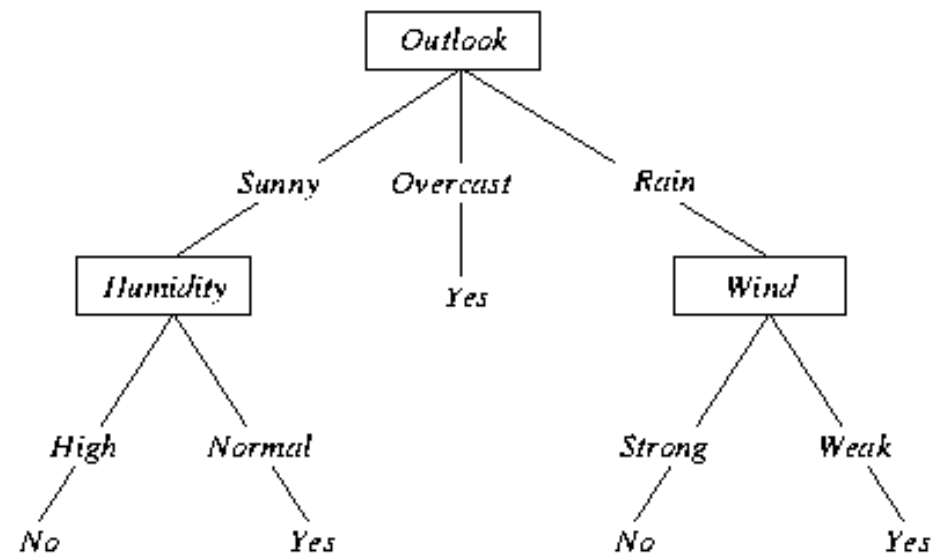
Day	Outlook	Temp	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Learned function $y(\mathbf{x})$: a tree



Components of Classification Tree

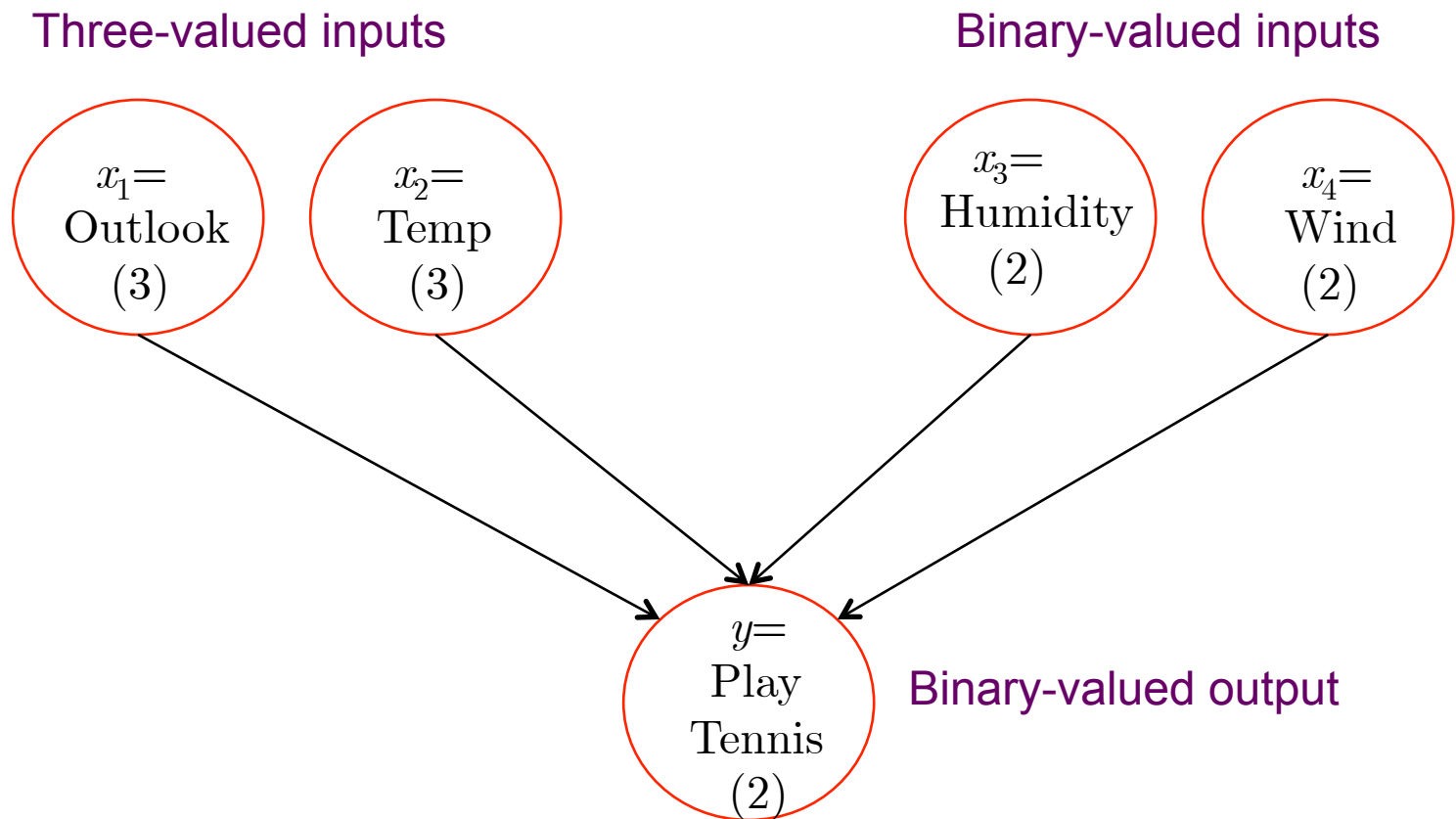
- Classify instances
 - by traversing from root to leaf node,
- Each node specifies a test of an attribute
- Each branch descending from a node corresponds a possible value of this attribute



Classification Tree Learning

- Tree approximates discrete-valued target function $y(\mathbf{x})$
 - Robust to noisy data and capable of learning disjunctive expressions
- A family of decision tree learning algorithms includes ID3, ASSISTANT and C4.5

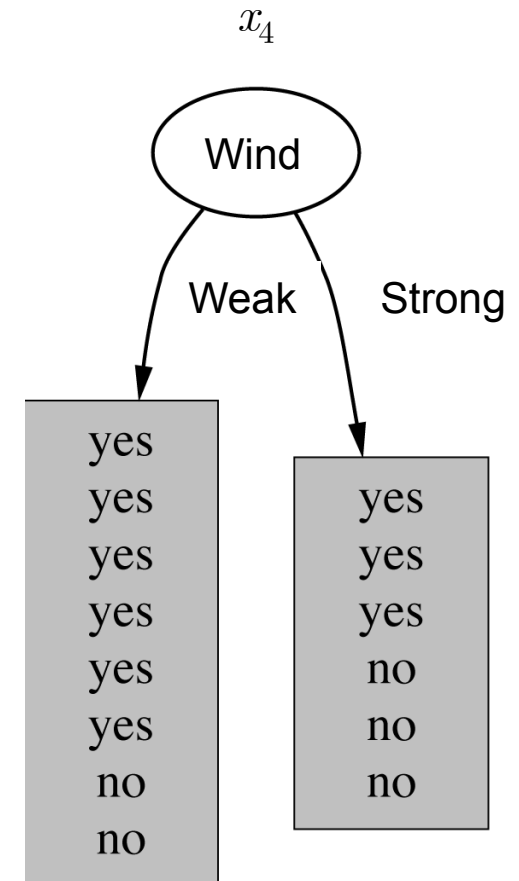
Graphical Model for PlayTennis



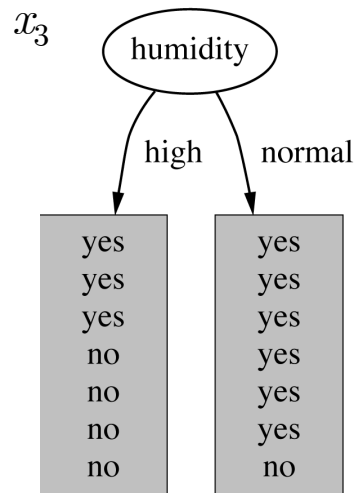
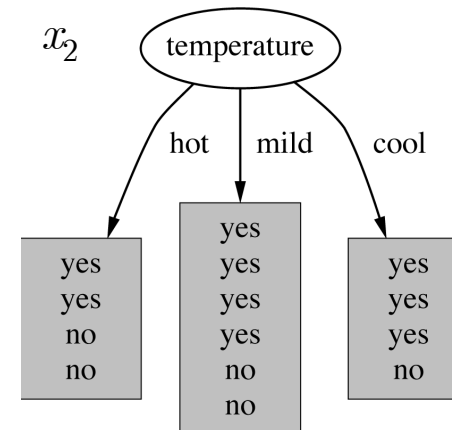
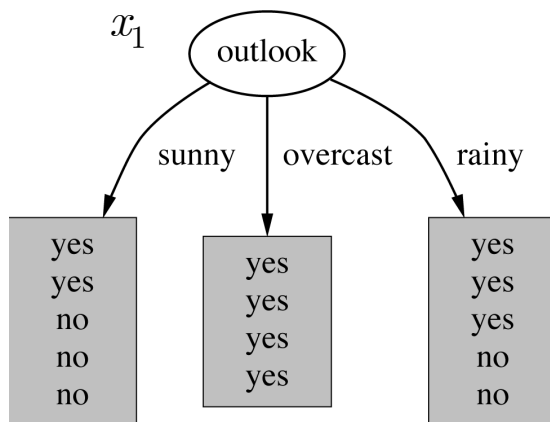
Need $3 \times 3 \times 2 \times 2 \times 2 = 76$ probabilities (parameters) for joint distribution $p(\mathbf{x}, y)$
Can we do better?

Tree Stump for a binary-valued input variable

Day	Outlook	Temp	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



Tree Stumps for other three inputs



Day	Outlook	Temp	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Good and Poor attributes

- An binary-valued attribute is good when:
 - for one value we get all instances as positive
 - for other value we get all instances as negative
- An attribute is poor when:
 - it provides no discrimination
 - attribute is immaterial to the decision
 - for each value we have same number of positive and negative instances

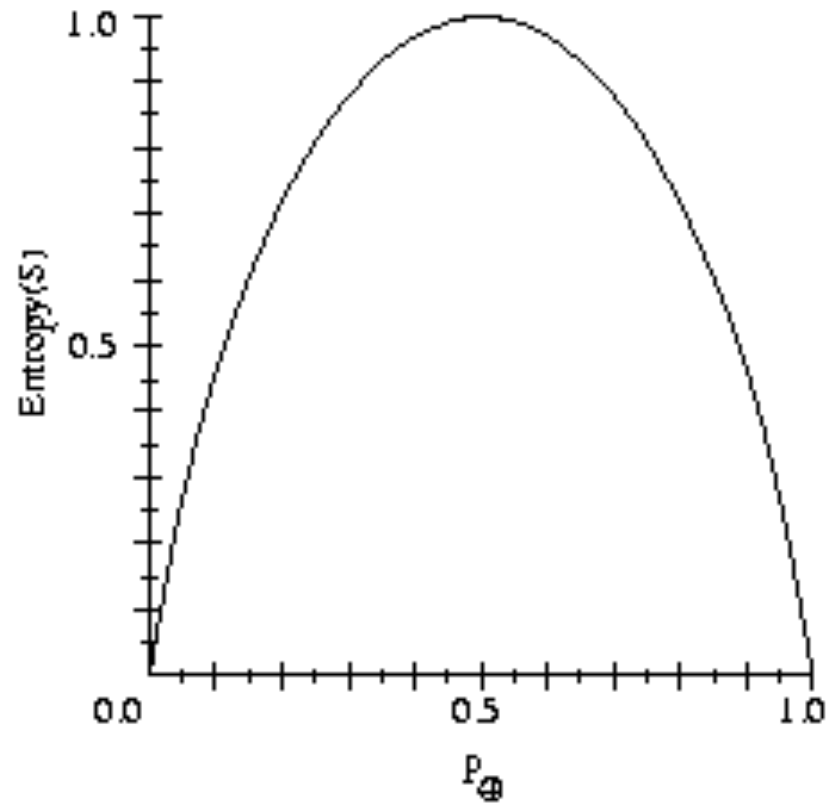
Entropy: Measure of Homogeneity of Examples

- Entropy: Characterizes the (im)purity of an arbitrary collection of examples
- Given a collection S of positive and negative examples, entropy of S relative to boolean classification is

$$\text{Entropy}(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_-$$

Where p_+ is proportion of positive examples and p_- is proportion of negative examples

Entropy Function Relative to a Boolean Classification



Entropy of data set

- Dataset S has 9 positive, 5 negative examples
- Entropy of S is:

$$\begin{aligned} \text{Entropy } (9+, 5-) = & \\ & -(9/14)\log_2(9/14) - \\ & (5/14)\log_2(5/14) = 0.94 \end{aligned}$$

Day	Outlook	Temp	Humidity	Wind	PlayTennis
D1	<i>Sunny</i>	<i>Hot</i>	<i>High</i>	<i>Weak</i>	<i>No</i>
D2	<i>Sunny</i>	<i>Hot</i>	<i>High</i>	<i>Strong</i>	<i>No</i>
D3	<i>Overcast</i>	<i>Hot</i>	<i>High</i>	<i>Weak</i>	<i>Yes</i>
D4	<i>Rain</i>	<i>Mild</i>	<i>High</i>	<i>Weak</i>	<i>Yes</i>
D5	<i>Rain</i>	<i>Cool</i>	<i>Normal</i>	<i>Weak</i>	<i>Yes</i>
D6	<i>Rain</i>	<i>Cool</i>	<i>Normal</i>	<i>Strong</i>	<i>No</i>
D7	<i>Overcast</i>	<i>Cool</i>	<i>Normal</i>	<i>Strong</i>	<i>Yes</i>
D8	<i>Sunny</i>	<i>Mild</i>	<i>High</i>	<i>Weak</i>	<i>No</i>
D9	<i>Sunny</i>	<i>Cool</i>	<i>Normal</i>	<i>Weak</i>	<i>Yes</i>
D10	<i>Rain</i>	<i>Mild</i>	<i>Normal</i>	<i>Weak</i>	<i>Yes</i>
D11	<i>Sunny</i>	<i>Mild</i>	<i>Normal</i>	<i>Strong</i>	<i>Yes</i>
D12	<i>Overcast</i>	<i>Mild</i>	<i>High</i>	<i>Strong</i>	<i>Yes</i>
D13	<i>Overcast</i>	<i>Hot</i>	<i>Normal</i>	<i>Weak</i>	<i>Yes</i>
D14	<i>Rain</i>	<i>Mild</i>	<i>High</i>	<i>Strong</i>	<i>No</i>

- Entropy is zero if all in S belong to the same class

Entropy for multi-valued target function

- If the target attribute can take on c different values, the entropy of S relative to this c -wise classification is

$$\text{Entropy} (S) \equiv \sum_{i=1}^c -p_i \log_2 p_i$$

Information Gain

- Entropy measures the impurity of a collection
- Information Gain is defined in terms of Entropy
 - expected reduction in entropy caused by partitioning the examples according to this attribute

Information Gain Measures the Expected Reduction in Entropy

- Information gain of attribute A is the reduction in entropy caused by partitioning the set of examples S

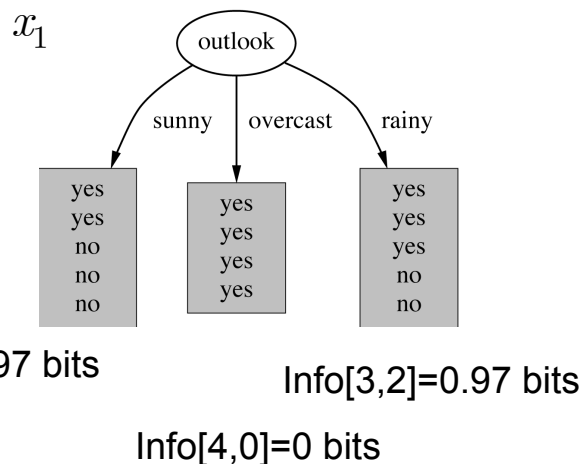
$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

- where $Values(A)$ is the set of all possible values for attribute A and S_v is the subset of S for which attribute A has value v

Information gain for attribute x_1

Outlook

For first value (*Sunny*) there are 2 positive and 3 negative examples



$$\begin{aligned} \text{Info}[2,3] &= \text{entropy}(2/5, 3/5) \\ &= -2/5 \log 2/5 - 3/5 \log 3/5 \\ &= 0.97 \text{ bits} \end{aligned}$$

Average info of subtree(weighted)=

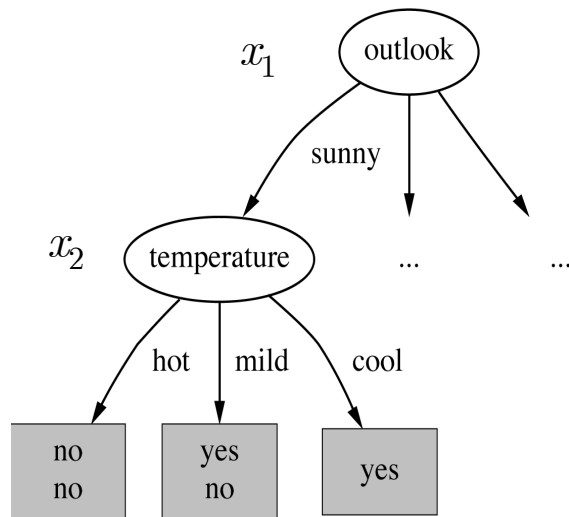
$$\begin{aligned} &0.97 \times 5/14 + \\ &0 \times 4/14 + \\ &0.97 \times 5/14 = 0.693 \text{ bits} \end{aligned}$$

Info of all training samples, $\text{info}[9,5] = 0.94$

$$\text{Gain}(S, \text{Outlook}) = 0.94 - 0.693 = 0.247 \text{ bits}$$

Expanded Tree Stumps for x_1

Outlook=Sunny



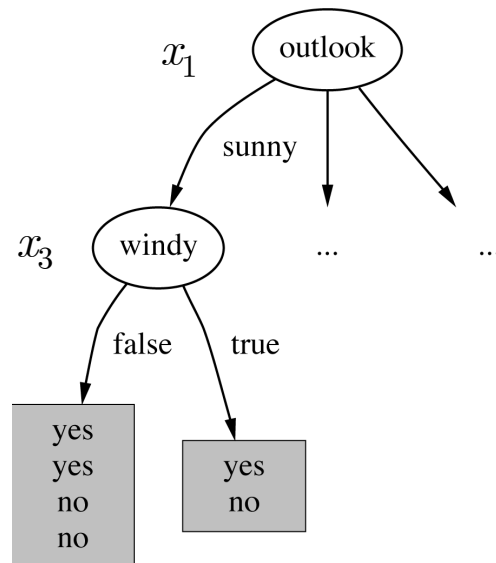
$$\text{Info}(2,3)=0.97$$

$$\text{Info}(0,2)=\text{info}(1,0)=0$$

$$\text{Info}(1,1)=0.5$$

$$\text{Ave Info}=0+0+(1/5)$$

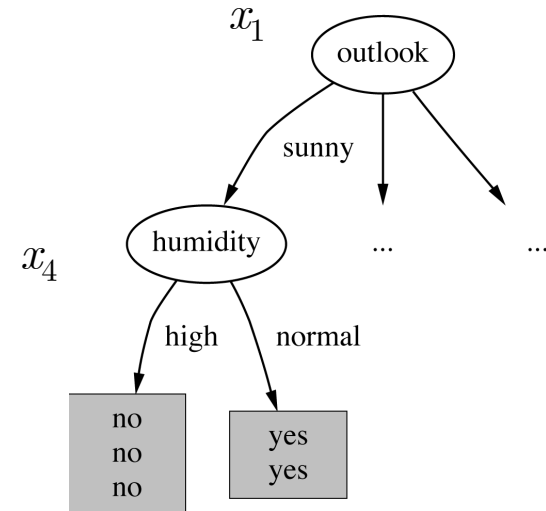
$$\text{Gain}(temp)=0.97-0.2=0.77 \text{ bits}$$



$$\text{Info}(3,3)=1$$

$$\text{Info}(2,2)=\text{Info}(1,1)=1$$

$$\text{Gain}(windy)=1-1=0 \text{ bits}$$



$$\text{Info}(3,2)=0.97$$

$$\text{Info}(0,3)=\text{Info}(2,0)=0$$

$$\text{Gain}(humidity)=0.97 \text{ bits}$$

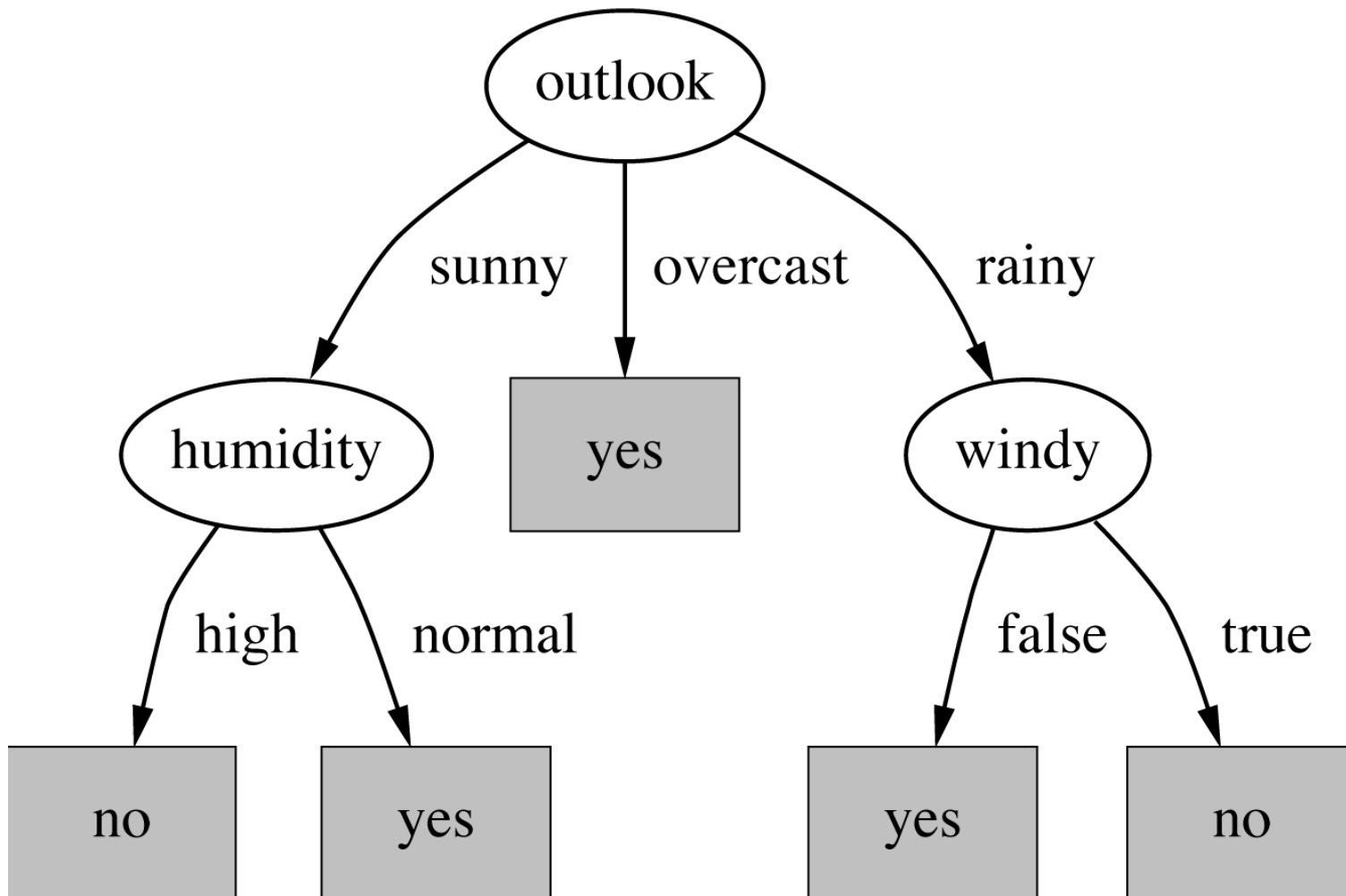
Since $\text{Gain}(humidity)$ is highest, select humidity as splitting attribute.

No need to split further

Information gain for each attribute

- $Gain(Outlook) = 0.94 - 0.693 = 0.247$
- $Gain(Temperature) = 0.94 - 0.911 = 0.029$
- $Gain(Humidity) = 0.94 - 0.788 = 0.152$
- $Gain(Windy) = 0.94 - 0.892 = 0.048$
- $\arg \max \{0.247, 0.029, 0.152, 0.048\} = Outlook$
- **Select *Outlook* as the splitting attribute of tree**

Decision Tree for the Weather Data



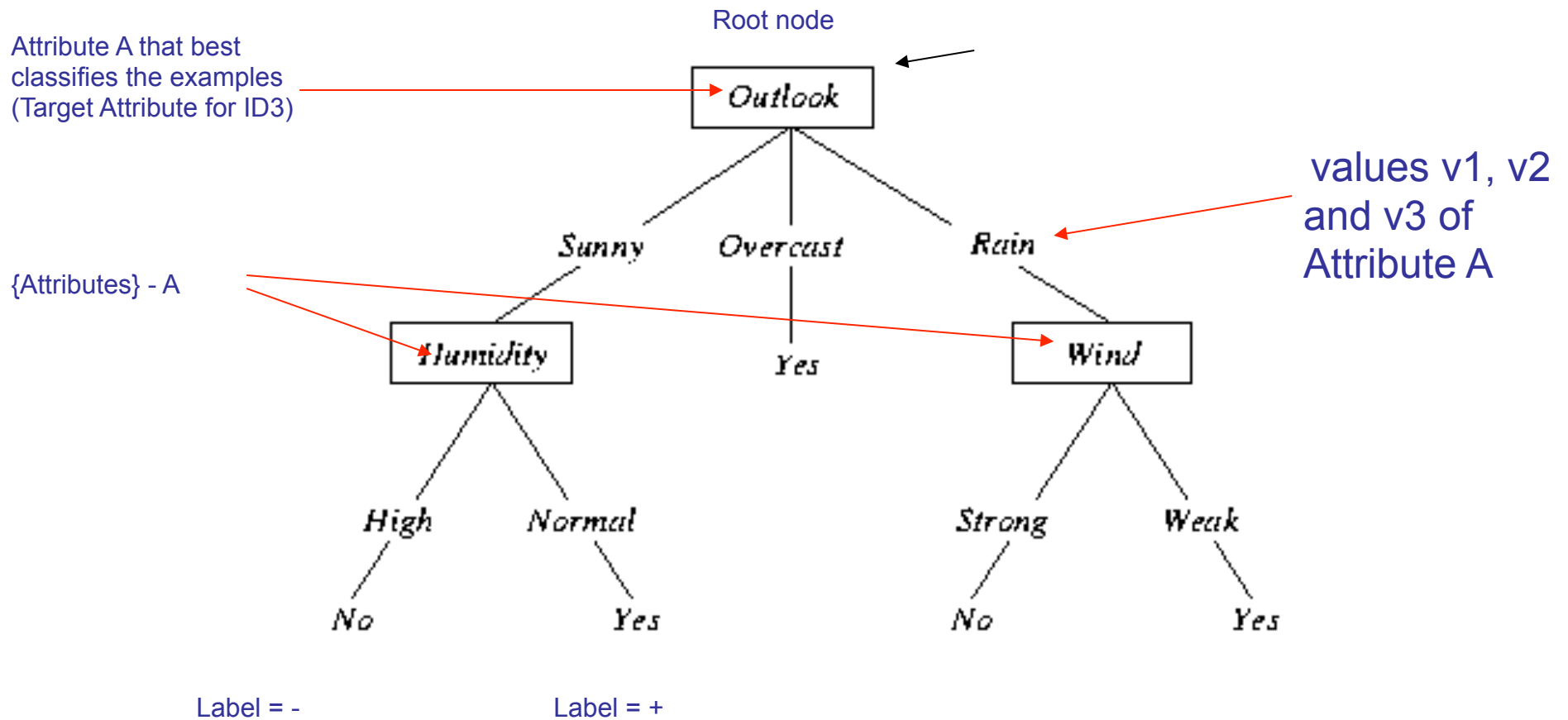
The Basic Decision Tree Learning Algorithm (ID3)

- Top-down, greedy search (no backtracking) through space of possible decision trees
- Begins with the question
 - “which attribute should be tested at the root of the tree?”
- Answer
 - evaluate each attribute to see how it alone classifies training examples
- Best attribute is used as root node
 - descendant of root node is created for each possible value of this attribute

Which Attribute Is Best for the Classifier?

- Select attribute that is most useful for classification
- ID3 uses Information gain as a quantitative measure of an attribute
- Information Gain:
 - A statistical property that measures how well a given attribute separates the training examples according to their target classification.

ID3 Algorithm Notation



ID3 Algorithm to learn boolean-valued functions

ID3 (Examples, Target_attribute, Attributes)

Examples are the training examples. Target_attribute is the attribute (or feature) whose value is to be predicted by the tree. Attributes is a list of other attributes that may be tested by the learned decision tree. Returns a decision tree (actually the root node of the tree) that correctly classifies the given Examples.

Create a *Root* node for the tree

- If all Examples are positive, Return the single-node tree *Root*, with label = +
- If all Examples are negative, Return the single-node tree *Root*, with label = -
- If Attributes is empty, Return the single-node tree *Root*, with label = the most common value of *Target_attribute* in *Examples*
- %Note that we will return the name of a feature at this point

ID3 Algorithm, continued

·Otherwise Begin

- A $\boxed{?}$ the attribute from *Attributes* that best* classifies *Examples*
- The decision attribute (**feature**) for *Root* $\boxed{?}$ A
- For each possible value v_j , of A,
 - Add a new tree branch below *Root*, corresponding to test $A = v_j$
 - Let $Examples_{v_j}$ the subset of *Examples* that have value v_j for A
 - If $Examples_{v_j}$ is empty'
 - Then below this new branch, add a leaf node with label = most common value of *Target_attribute* in *Examples*
 - Else, below this new branch add the subtree

$$ID3(Examples_{v_j}, Target_attribute, Attributes - \{A\})$$

·End

·Return *Root*

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

- The best attribute is the one with the highest *information gain*, as defined in Equation:

Perfect feature

- If feature outlook has two values: sunny and rainy
- If for sunny all 5 values of playtennis are yes
- If for rainy all 9 values of playtennis are no

$$\begin{aligned}\text{Gain}(S, \text{outlook}) &= 0.94 - (5/9).0 - (9/9).0 \\ &= 0.94\end{aligned}$$

Training Examples for Target Concept *PlayTennis*

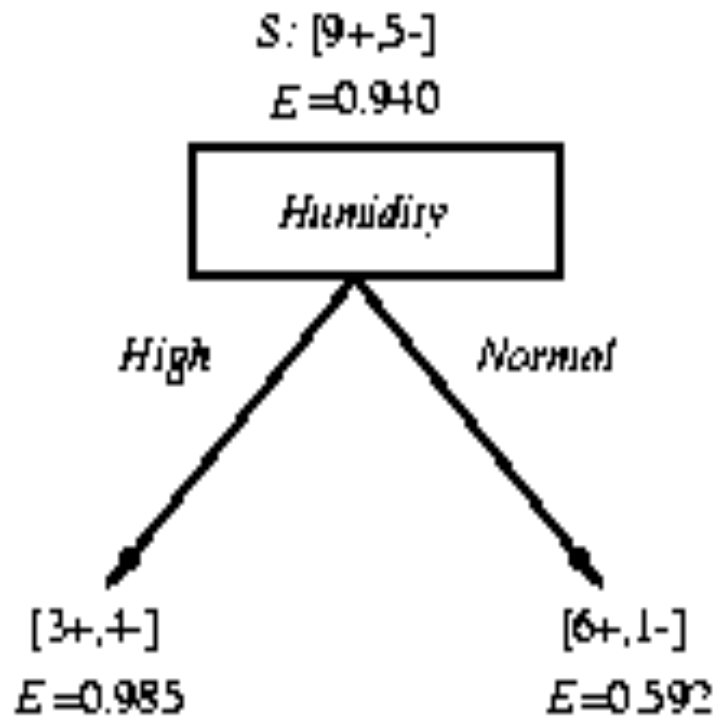
Day	Outlook	Temp	Humidity	Wind	<i>PlayTennis</i>
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Table 3.2

Stepping through ID3 for the example

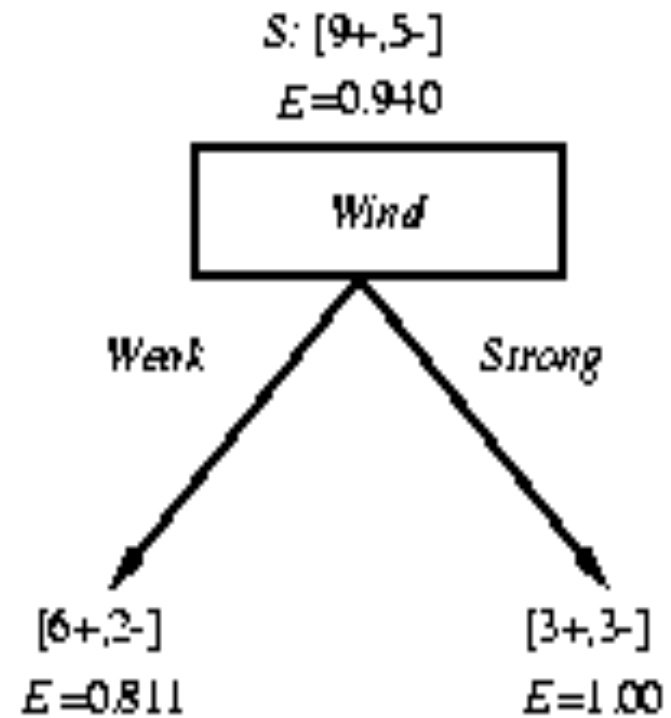
- **Top Level (with $S = \{D1, \dots, D14\}$)**
 - $\text{Gain}(S, \text{Outlook}) = 0.246$
 - $\text{Gain}(S, \text{Humidity}) = 0.151$ ← Best prediction of target attribute
 - $\text{Gain}(S, \text{Wind}) = 0.048$
 - $\text{Gain}(S, \text{Temperature}) = 0.029$
- **Example computation of Gain for Wind**
 - $\text{Values}(\text{Wind}) = \text{Weak}, \text{Strong}$
 - $S = [9+, 5-]$
 - $S_{\text{weak}} \leftarrow [6+, 2-], S_{\text{strong}} \leftarrow [3+, 3-]$
 - $\text{Gain}(S, \text{Wind}) = 0.940 - (8/14)0.811 - (6/14)1.00$
 $= 0.048$

Sample calculations for Humidity and Wind



$Gain(S, Humidity)$

$$\begin{aligned}
 &= 0.940 - (7/14) \cdot 0.985 - (7/14) \cdot 0.592 \\
 &= .151
 \end{aligned}$$

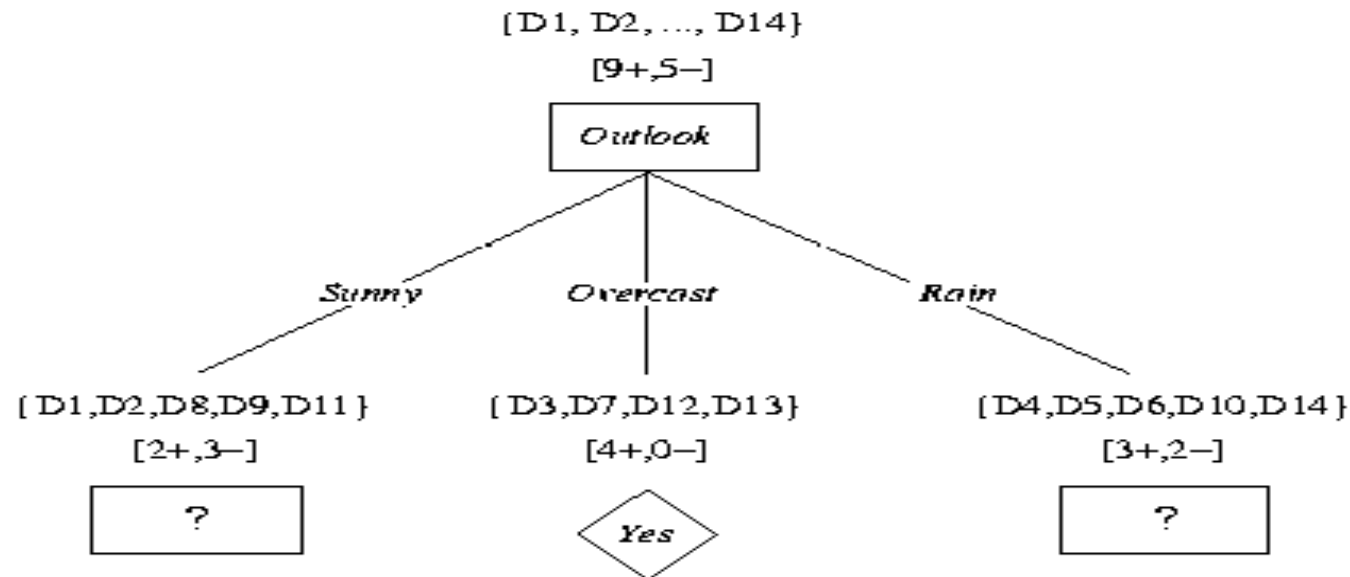


$Gain(S, Wind)$

$$\begin{aligned}
 &= 0.940 - (8/14) \cdot 0.811 - (6/14) \cdot 1.0 \\
 &= .048
 \end{aligned}$$

Figure 3.3

The Partially Learned Decision Tree



Which attribute should be tested here?

$$S_{\text{sunny}} = \{D1, D2, D8, D9, D11\}$$

$$\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = .970 - (3/5)0.0 - (2/5)0.0 = .970$$

$$\text{Gain}(S_{\text{sunny}}, \text{Temperature}) = .970 - (2/5)0.0 - (2/5)1.0 - (1/5)0.0 = .570$$

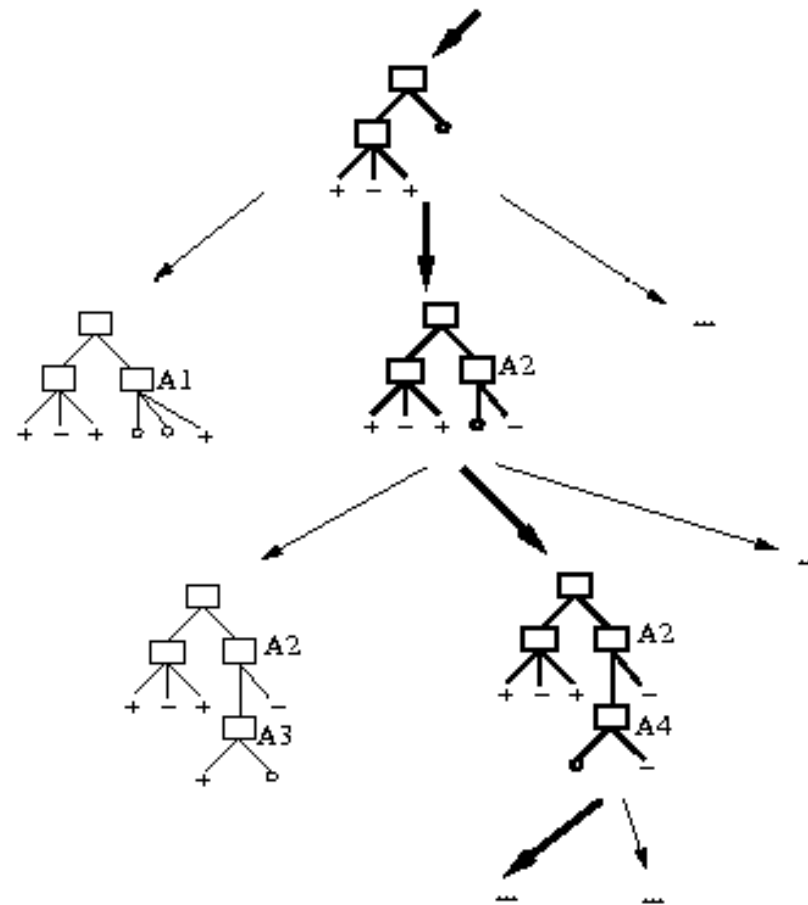
$$\text{Gain}(S_{\text{sunny}}, \text{Wind}) = .970 - (2/5)1.0 - (3/5).918 = .019$$

Humidity is a perfect feature since there is no uncertainty left when we know its value

Hypothesis Space Search in Decision Tree Learning

- ID3 searches a hypothesis space for one that fits training examples
- Hypothesis space searched is set of possible decision trees
- ID3 performs hill-climbing, starting with empty tree, considering progressively more elaborate hypotheses (to find tree to correctly classify training data)
- Hill climbing is guided by evaluation function which is the gain measure

Hypothesis Space Search by ID3

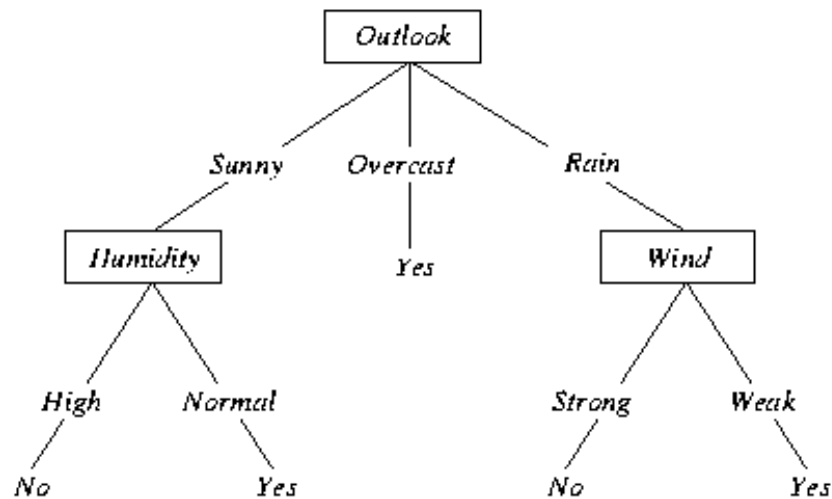


Information gain
heuristic
guides search of
hypothesis space by ID3

Figure 3.5

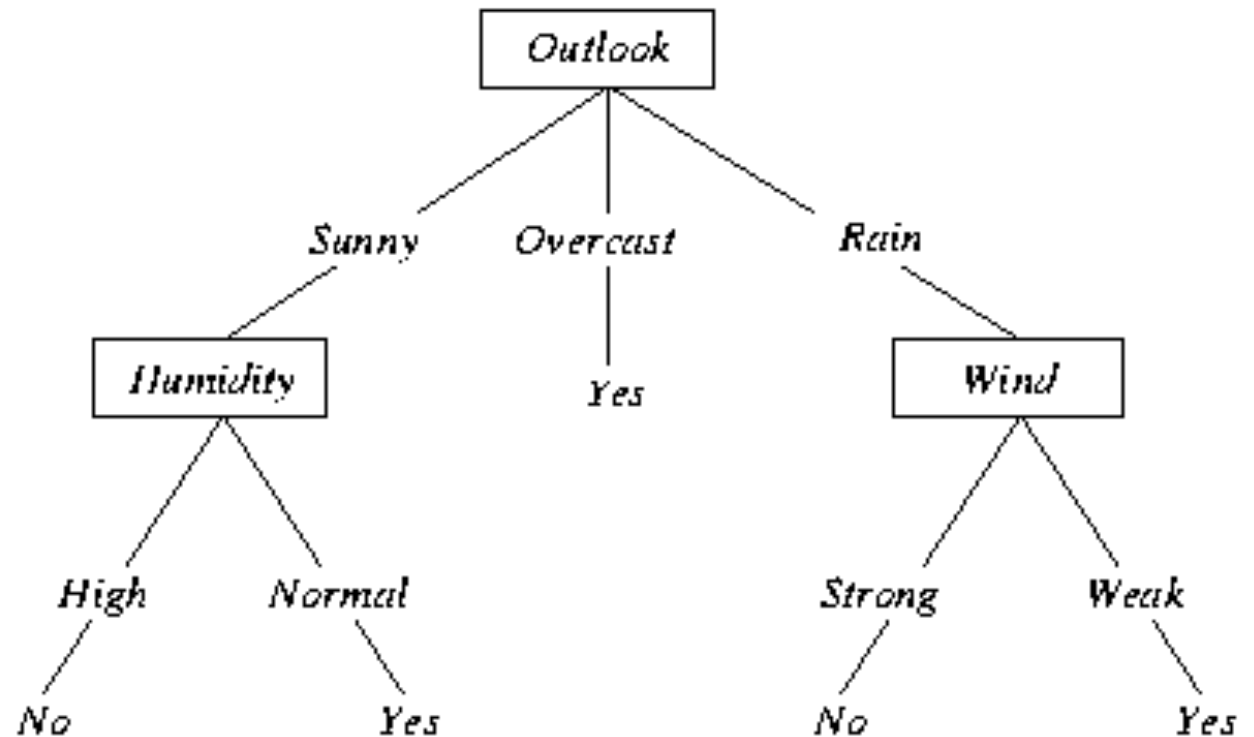
Decision Trees and Rule-Based Systems

- Learned trees can also be re-represented as sets of *if-then* rules to improve human readability



Decision Trees represent disjunction of conjunctions

$(\text{Outlook} = \text{Sunny} \wedge \text{Humidity} = \text{Normal})$
 $\vee (\text{Outlook} = \text{Overcast})$
 $\vee (\text{Outlook} = \text{Rain} \wedge \text{Wind} = \text{Weak})$



Appropriate Problems for Decision Tree Learning

- Instances are represented by attribute-value pairs
 - each attribute takes on a small no of disjoint possible values, eg, hot, mild, cold
 - extensions allow real-valued variables as well, eg temperature
- The target function has discrete output values
 - eg, Boolean classification (yes or no)
 - easily extended to multiple-valued functions
 - can be extended to real-valued outputs as well

Appropriate Problems for Decision Tree Learning (2)

- Disjunctive descriptions may be required
 - naturally represent disjunctive expressions
- The training data may contain errors
 - robust to errors in classifications and in attribute values
- The training data may contain missing attribute values
 - eg, humidity value is known only for some training examples

Appropriate Problems for Decision Tree Learning (3)

- Practical problems that fit these characteristics are:
 - learning to classify
 - medical patients by their disease
 - equipment malfunctions by their cause
 - loan applications by by likelihood of defaults on payments

Capabilities and Limitations of ID3

- Hypothesis space is a complete space of all discrete valued functions
- Cannot determine how many alternative trees are consistent with training data (follows from maintaining a single current hypothesis)
- ID3 in its pure form performs no backtracking (usual risks of hill-climbing- converges to local optimum)
- ID3 uses all training examples at each step to make statistically based decisions regarding how to refine its current hypothesis
 - more robust than Find-S and Candidate Elimination which are incrementally-based