# Decision Trees

These slides were assembled by Eric Eaton, with grateful acknowledgement of the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution. Please send comments and corrections to Eric.

# Function Approximation

**Problem Setting**

- Set of possible instances $\mathcal{X}$

- Set of possible labels $\mathcal{Y}$

- Unknown target function $f : \mathcal{X} \to \mathcal{Y}$

- Set of function hypotheses $H = \{h \mid h : \mathcal{X} \to \mathcal{Y}\}$

**Input**: Training examples of unknown target function $f$

$$\{\langle \boldsymbol{x}_i, y_i \rangle\}_{i=1}^{n} = \{\langle \boldsymbol{x}_1, y_1 \rangle, \dots, \langle \boldsymbol{x}_n, y_n \rangle\}$$

**Output**: Hypothesis $h \in H$ that best approximates $f$
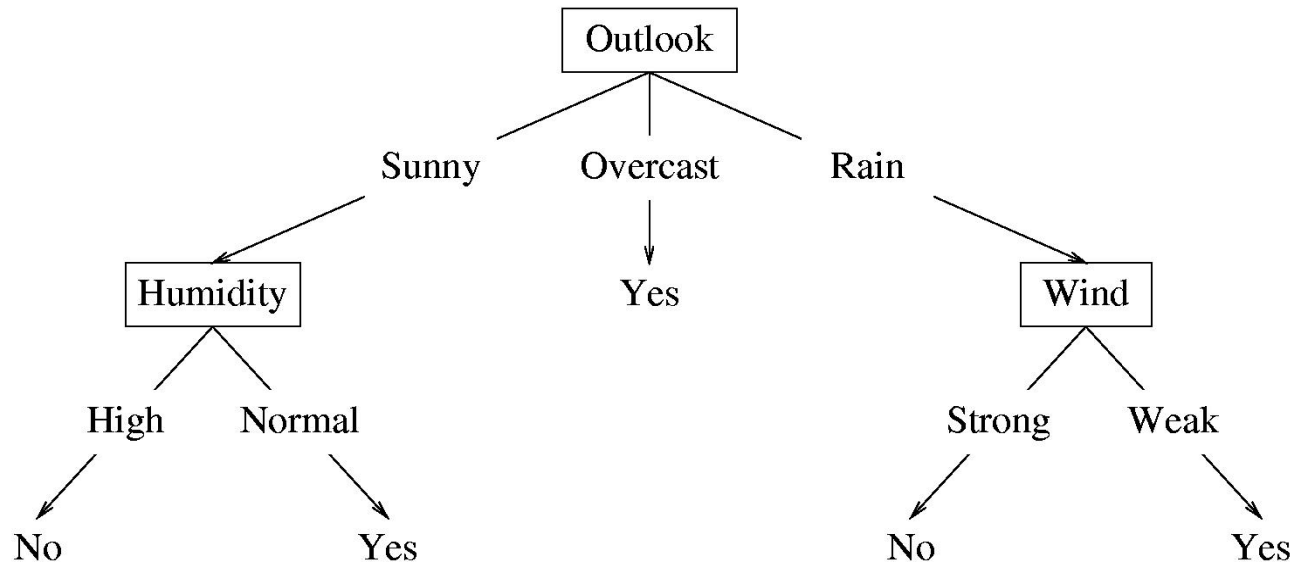
# Sample Dataset

- Columns denote features $X_i$
- Rows denote labeled instances $\langle \boldsymbol{x}_i, y_i \rangle$
- Class label denotes whether a tennis game was played

$\langle \boldsymbol{x}_i, y_i \rangle$

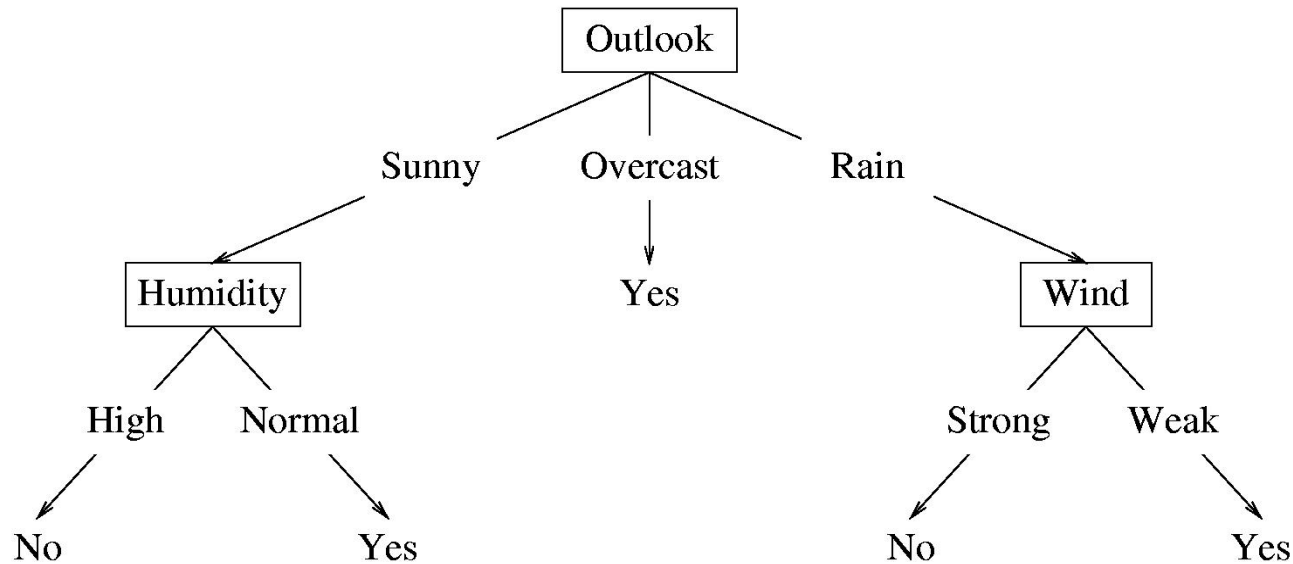| | Predictors | | | Response |
|---|---|---|---|---|
| **Outlook** | **Temperature** | **Humidity** | **Wind** | **Class** |
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

# Decision Tree

- A possible decision tree for the data:



- Each internal node: test one attribute $X_i$
- Each branch from a node: selects one value for $X_i$
- Each leaf node: predict $Y$ (or $p(Y \mid \boldsymbol{x} \in \mathrm{leaf})$ )

# Decision Tree

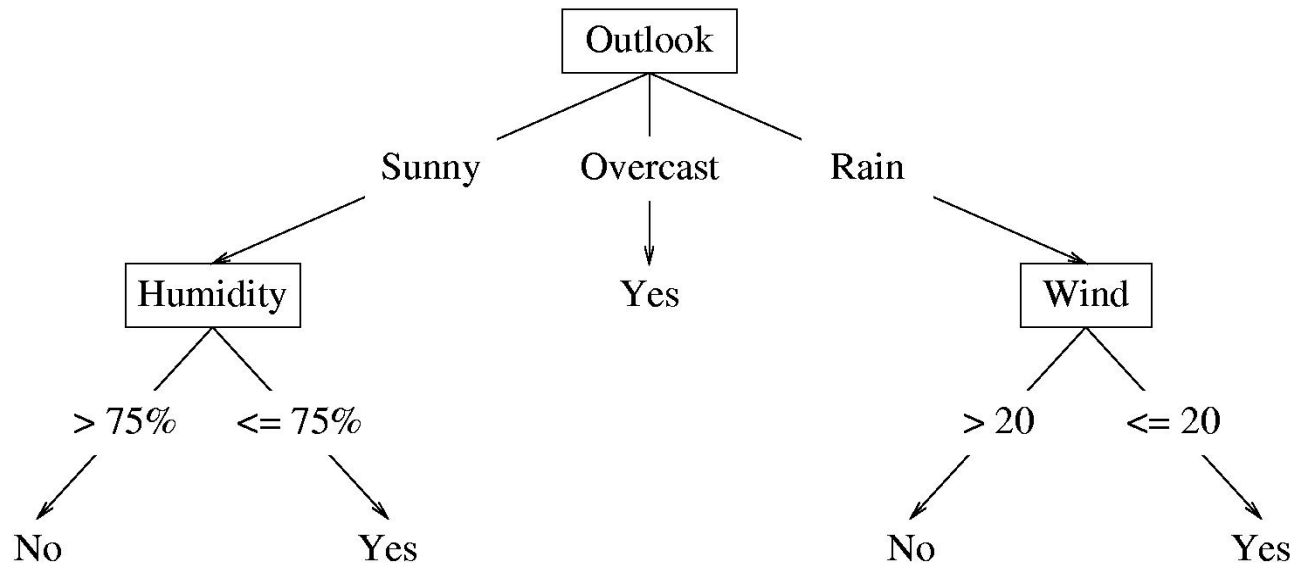- A possible decision tree for the data:



- What prediction would we make for

<outlook=sunny, temperature=hot, humidity=high, wind=weak> ?
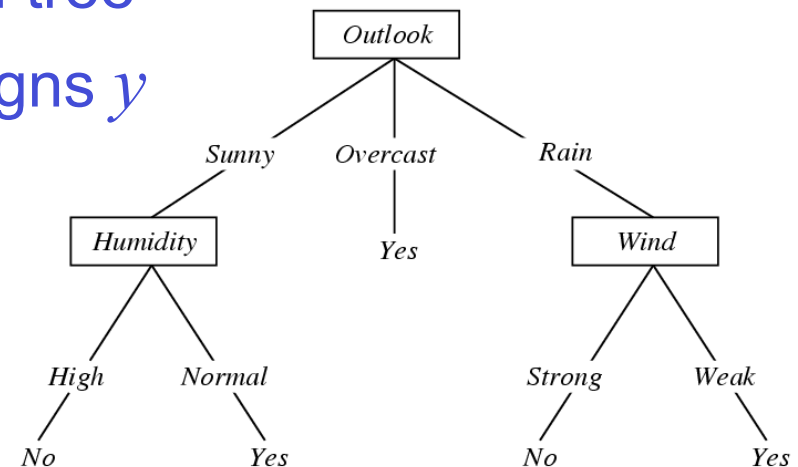
# Decision Tree

- If features are continuous, internal nodes can test the value of a feature against a threshold

# Decision Tree Learning

**Problem Setting**:

- Set of possible instances $X$

  – each instance $x$ in $X$ is a feature vector

  – e.g., *<Humidity=low, Wind=weak, Outlook=rain, Temp=hot>*

- Unknown target function $f : X \rightarrow Y$

  – *Y* is discrete valued

- Set of function hypotheses $H=\{\, h \mid h : X \rightarrow Y \,\}$

  – each hypothesis $h$ is a decision tree
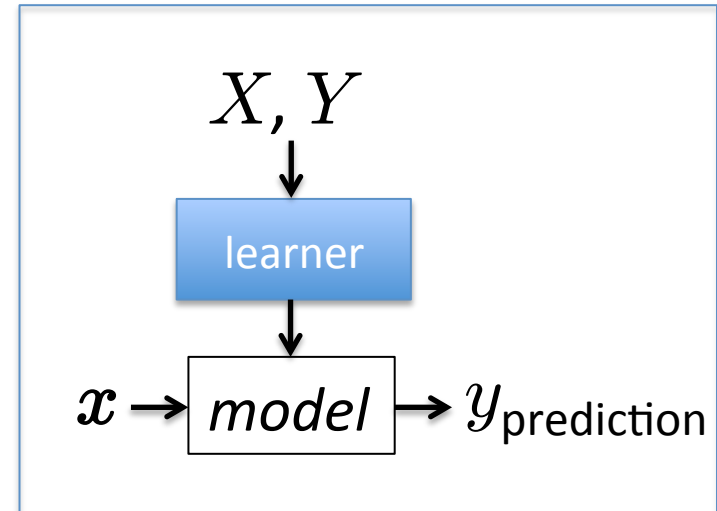
  – trees sorts $x$ to leaf, which assigns $y$

# Stages of (Batch) Machine Learning

**Given:** labeled training data $X, Y = \{\langle \boldsymbol{x}_i, y_i \rangle\}_{i=1}^{n}$

- Assumes each $\boldsymbol{x}_i \sim \mathcal{D}(\mathcal{X})$ with $y_i = f_{target}(\boldsymbol{x}_i)$

**Train the model:**

$model \leftarrow classifier.\text{train}(X, Y)$

$$
\begin{array}{c}
X, Y \\
\downarrow \\
\boxed{\text{learner}} \\
\downarrow \\
\boldsymbol{x} \rightarrow \boxed{model} \rightarrow y_{\text{prediction}}
\end{array}
$$

**Apply the model to new data:**

- Given: new unlabeled instance $\boldsymbol{x} \sim \mathcal{D}(\mathcal{X})$

$y_{\text{prediction}} \leftarrow model.\text{predict}(\boldsymbol{x})$

# Basic Algorithm for Top-Down Induction of Decision Trees

[ID3, C4.5 by Quinlan]

*node* = root of decision tree

Main loop:

1.  *A* ← the "best" decision attribute for the next node.
2.  Assign *A* as decision attribute for *node*.
3.  For each value of *A*, create a new descendant of *node*.
4.  Sort training examples to leaf nodes.
5.  If training examples are perfectly classified, stop. Else, recurse over new leaf nodes.
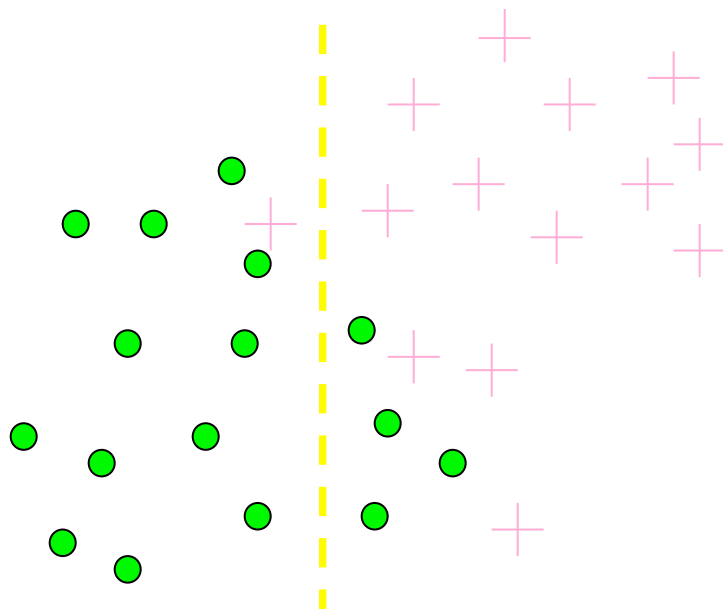
How do we choose which attribute is best?

# Choosing the Best Attribute

**Key problem**: choosing which attribute to split a given set of examples

- Some possibilities are:
  - **Random:** Select any attribute at random
  - **Least-Values:** Choose the attribute with the smallest number of possible values
  - **Most-Values:** Choose the attribute with the largest number of possible values
  - **Max-Gain:** Choose the attribute that has the largest expected *information gain*
    - i.e., attribute that results in smallest expected size of subtrees rooted at its children

- The ID3 algorithm uses the Max-Gain method of selecting the best attribute
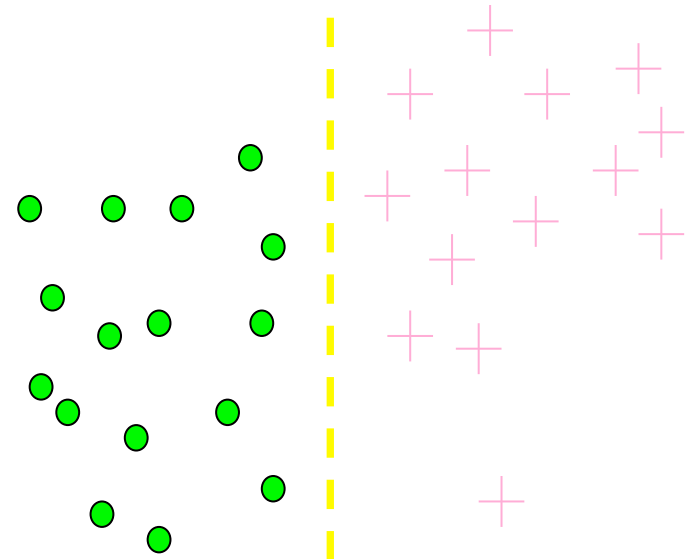
# Information Gain

Which test is more informative?

**Split over whether Balance exceeds 50K**

Less or equal 50K    Over 50K
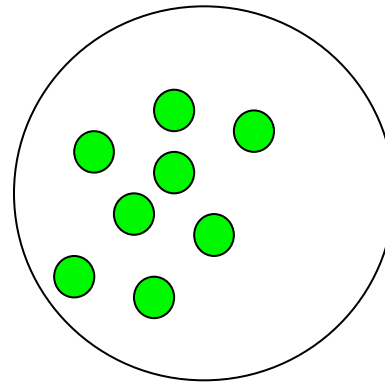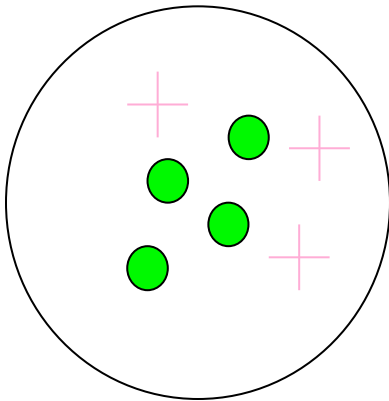
**Split over whether applicant is employed**
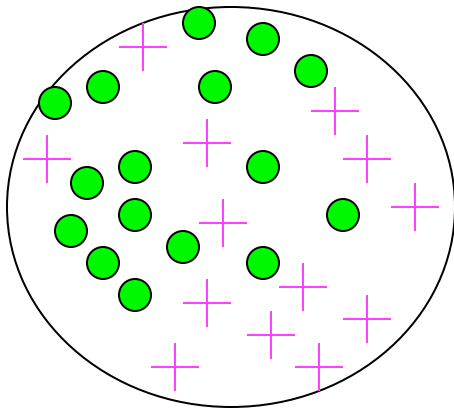
Unemployed    Employed

# Information Gain

## **Impurity/Entropy** (informal)

 – Measures the level of **impurity** in a group of examples

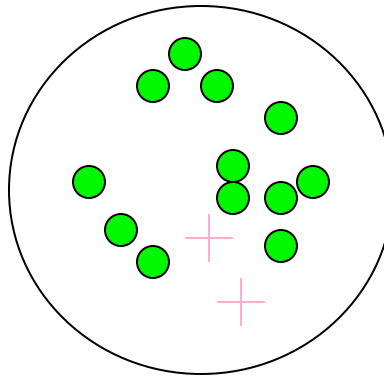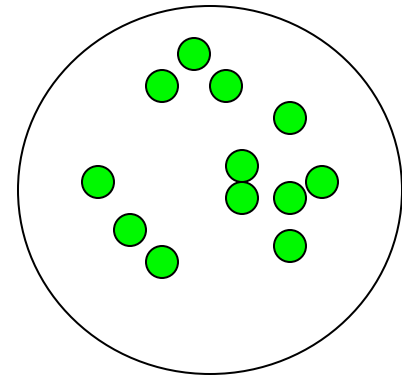# Impurity

**Very impure group**

**Less impure**

**Minimum impurity**

Based on slide by Pedro Domingos

# 2-Class Cases:
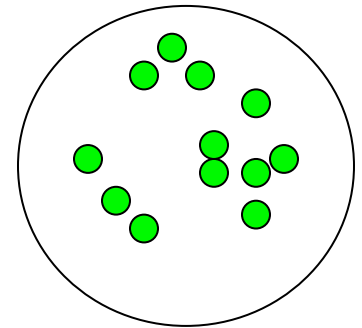
Entropy $$H(x) = -\sum_{i=1}^{n} P(x=i) \log_2 P(x=i)$$

- What is the entropy of a group in which all examples belong to the same class?
  - entropy = - 1 $\log_2$1 = 0

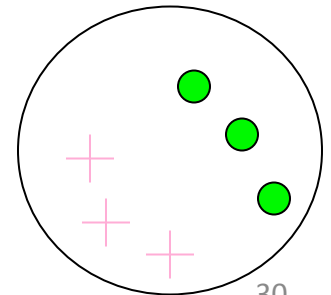  not a good training set for learning

**Minimum impurity**



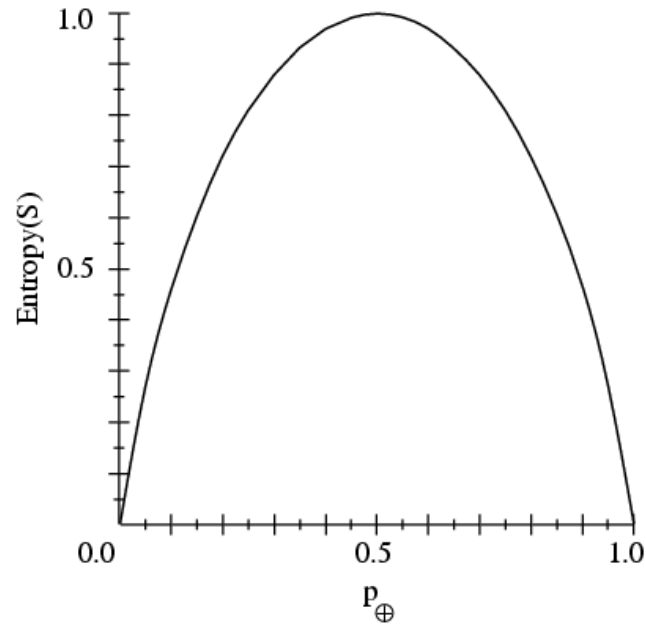- What is the entropy of a group with 50% in either class?
  - entropy = -0.5 $\log_2$0.5 − 0.5 $\log_2$0.5 =1

  good training set for learning

**Maximum impurity**



30

Based on slide by Pedro Domingos

# Sample Entropy



- $S$ is a sample of training examples
- $p_\oplus$ is the proportion of positive examples in $S$
- $p_\ominus$ is the proportion of negative examples in $S$
- Entropy measures the impurity of $S$

$$H(S) \equiv -p_\oplus \log_2 p_\oplus - p_\ominus \log_2 p_\ominus$$

# Information Gain

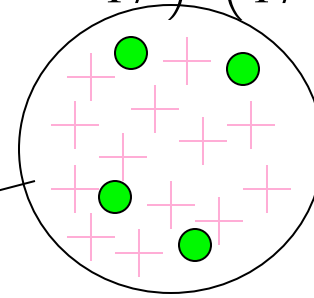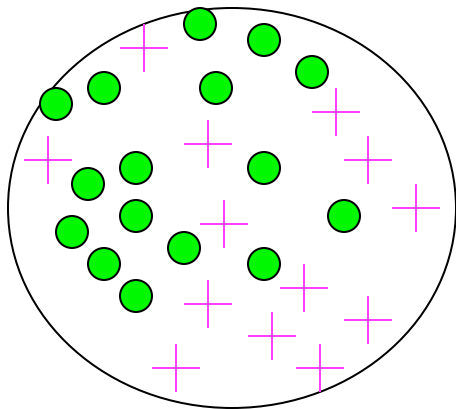- We want to determine which attribute in a given set of training feature vectors is most useful for discriminating between the classes to be learned.

- Information gain tells us how important a given attribute of the feature vectors is.

- We will use it to decide the ordering of attributes in the nodes of a decision tree.

# Calculating Information Gain

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

child
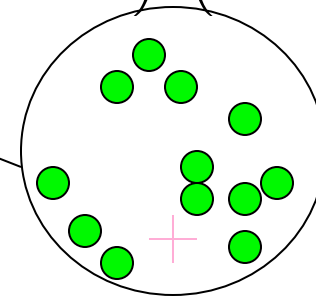entropy $-\left(\dfrac{13}{17}\cdot\log_2\dfrac{13}{17}\right)-\left(\dfrac{4}{17}\cdot\log_2\dfrac{4}{17}\right)=0.787$

Entire population (30 instances)

17 instances

child
entropy $-\left(\dfrac{1}{13}\cdot\log_2\dfrac{1}{13}\right)-\left(\dfrac{12}{13}\cdot\log_2\dfrac{12}{13}\right)=0.391$

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

13 instances

(Weighted) Average Entropy of Children = $\left(\dfrac{17}{30}\cdot 0.787\right)+\left(\dfrac{13}{30}\cdot 0.391\right)=0.615$
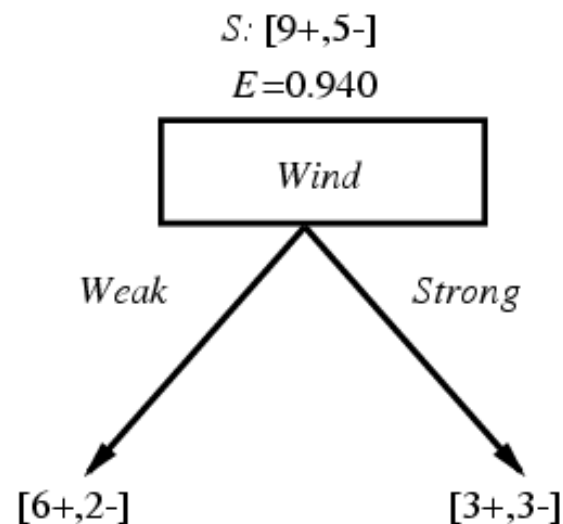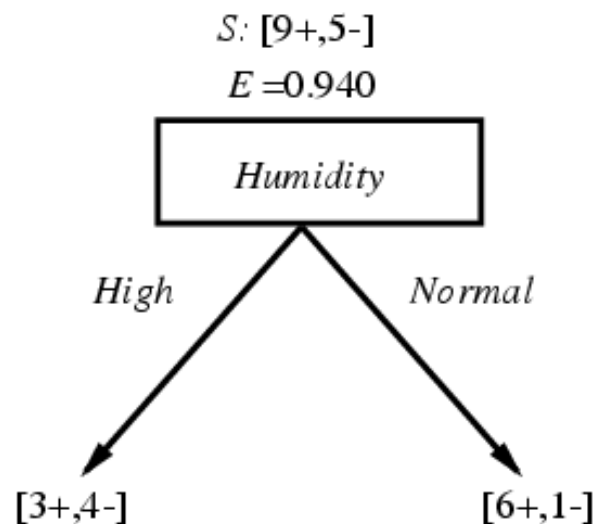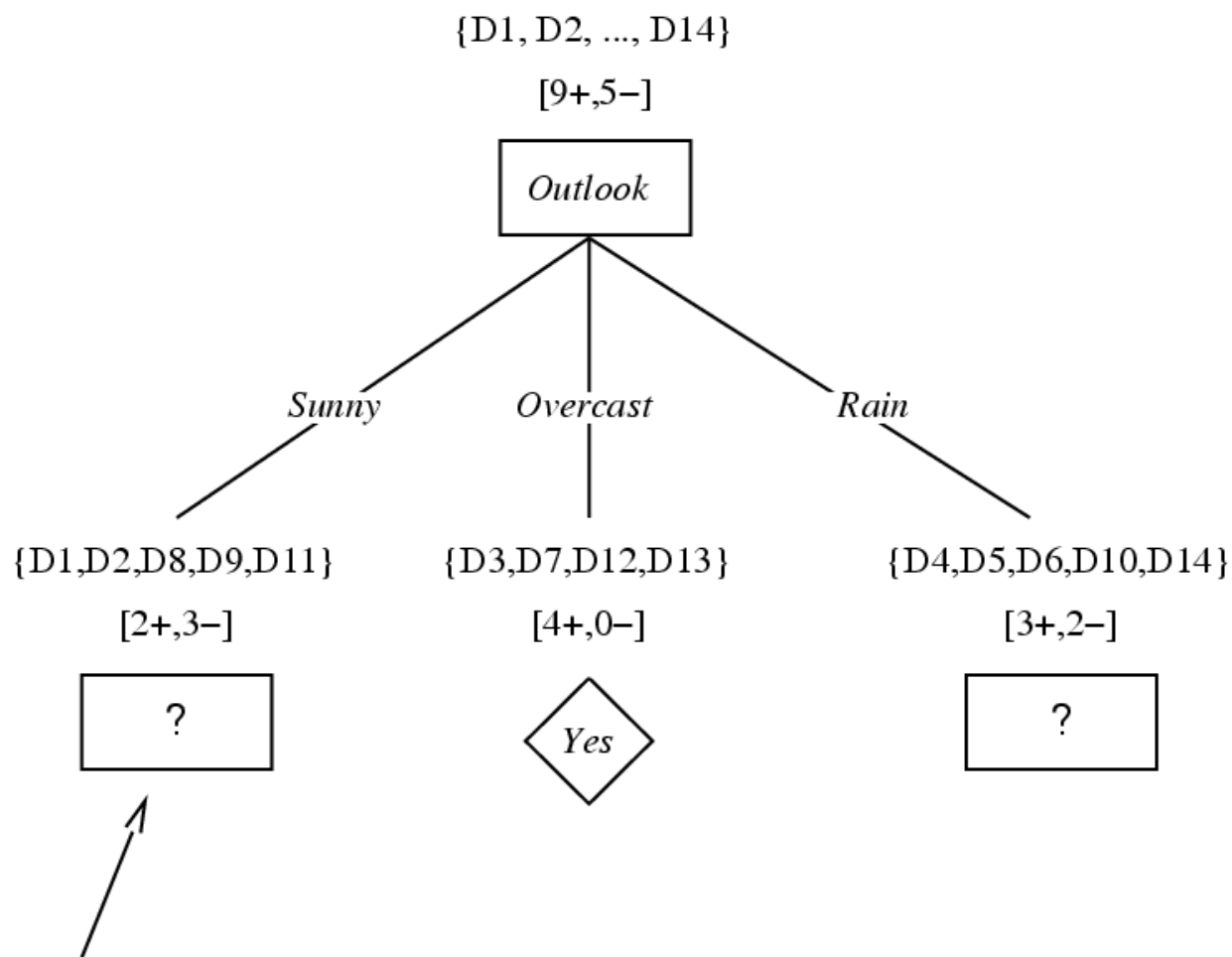
**Information Gain= 0.996 - 0.615 = 0.38**

# Training Examples

| Day | Outlook | Temperature | 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 |

# Selecting the Next Attribute

**Which attribute is the best classifier?**



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

Humidity

High                    Normal

[3+,4-]                 [6+,1-]

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

Wind

Weak                    Strong

[6+,2-]                 [3+,3-]

{D1, D2, ..., D14}

[9+,5−]

Outlook

Sunny          Overcast          Rain

{D1,D2,D8,D9,D11}      {D3,D7,D12,D13}      {D4,D5,D6,D10,D14}

[2+,3−]          [4+,0−]          [3+,2−]

?          Yes          ?

*Which attribute should be tested here?*

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

*Gain* ($S_{sunny}$ , *Humidity*) = .970 − (3/5) 0.0 − (2/5) 0.0 = .970

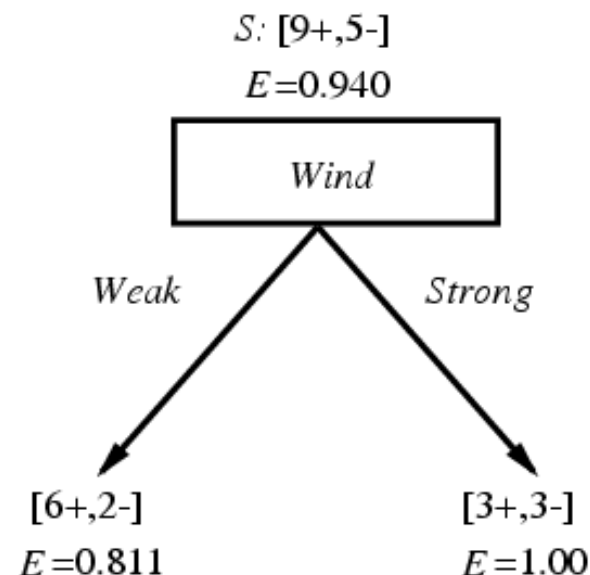*Gain* ($S_{sunny}$ , *Temperature*) = .970 − (2/5) 0.0 − (2/5) 1.0 − (1/5) 0.0 = .570

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

# Selecting the Next Attribute

**Which attribute is the best classifier?**

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

Humidity

High         Normal

[3+,4-]
*E* =0.985

[6+,1-]
*E* =0.592

*Gain (S, Humidity )*

= .940 - (7/14).985 - (7/14).592
= .151

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

Wind

Weak        Strong

[6+,2-]
*E* =0.811

[3+,3-]
*E* =1.00

*Gain (S, Wind)*

= .940 - (8/14).811 - (6/14)1.0
= .048