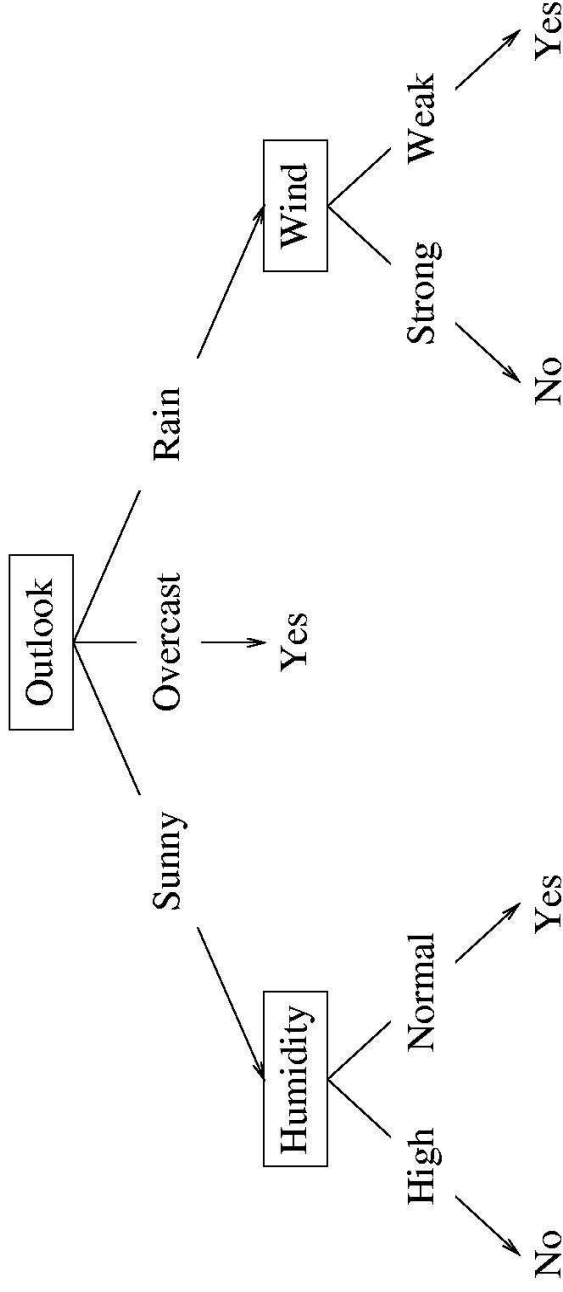


# Decision Trees

## Decision Tree Hypothesis Space

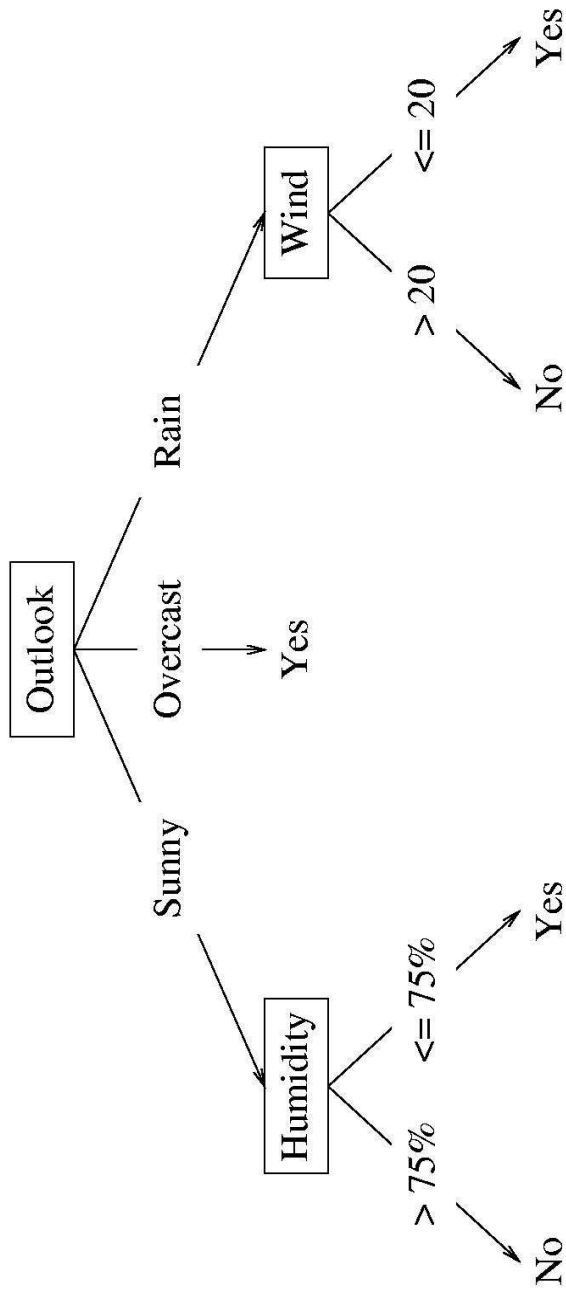
- **Internal nodes** test the value of particular features  $x_j$  and branch according to the results of the test.
- **Leaf nodes** specify the class  $h(\mathbf{x})$ .



Suppose the features are **Outlook** ( $x_1$ ), **Temperature** ( $x_2$ ), **Humidity** ( $x_3$ ), and **Wind** ( $x_4$ ). Then the feature vector  $\mathbf{x} = (\text{Sunny}, \text{Hot}, \text{High}, \text{Strong})$  will be classified as **No**. The **Temperature** feature is irrelevant.

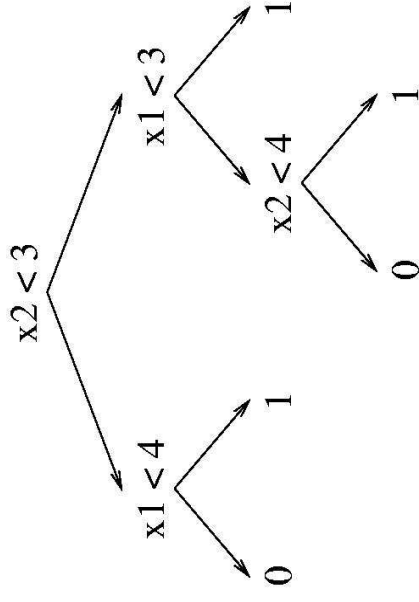
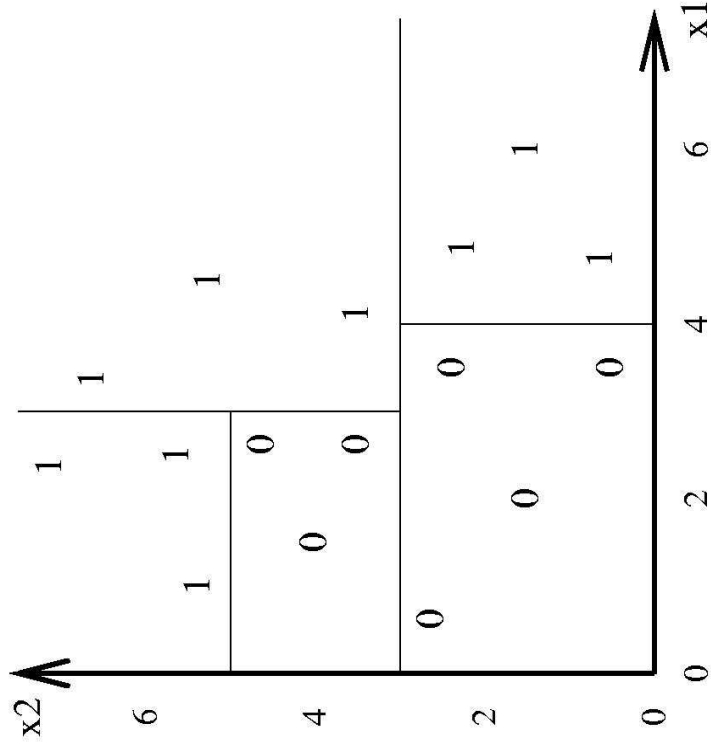
## Decision Tree Hypothesis Space

If the features are continuous, internal nodes may test the value of a feature against a threshold.

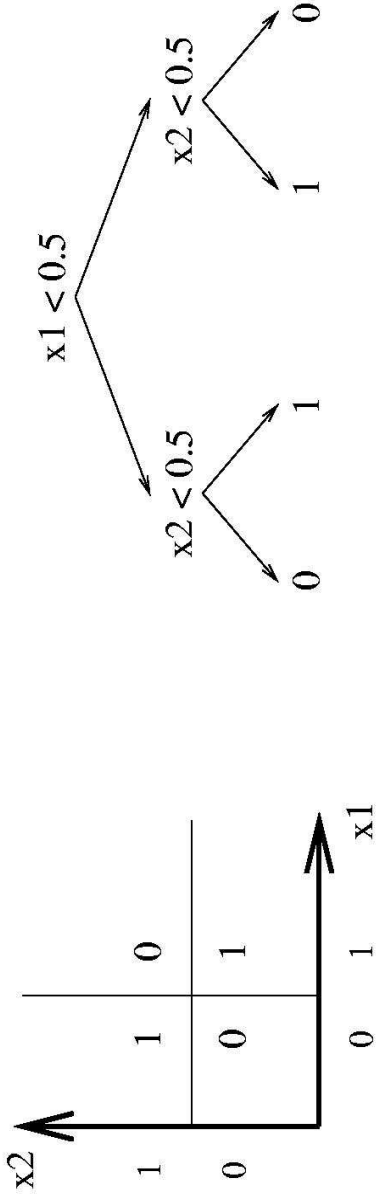


## Decision Tree Decision Boundaries

Decision trees divide the feature space into axis-parallel rectangles, and label each rectangle with one of the  $K$  classes.



## Decision Trees Can Represent Any Boolean Function



The tree will in the worst case require exponentially many nodes, however.

# The “No Free Lunch” Theorem

- Is there **any** representation that is compact (ie, sub-exponential in  $n$ ) for **all** functions?
  - Function = truth table
  - $n$  attributes  $\square 2^n$  rows in table
  - Classification/target column is  $2^n$  long
  - If you drop a bit, you cut the number of functions in half!
- 6 attributes = 18,446,744,073,709,551,616 functions

## Decision Trees Provide Variable-Size Hypothesis Space

As the number of nodes (or depth) of tree increases, the hypothesis space grows

- **depth 1** (“decision stump”) can represent any boolean function of one feature.
- **depth 2** Any boolean function of two features; some boolean functions involving three features (e.g.,  $(x_1 \wedge x_2) \vee (\neg x_1 \wedge \neg x_3)$ )
- **etc.**

## Learning Decision Trees

Decision trees provide a very popular and efficient hypothesis space.

- **Variable Size.** Any boolean function can be represented.
- **Deterministic.**
- **Discrete and Continuous Parameters.**

Learning algorithms for decision trees can be described as

- **Constructive Search.** The tree is built by adding nodes.
- **Eager.**
- **Batch** (although online algorithms do exist).



## Learning Algorithm for Decision Trees

The same basic learning algorithm has been discovered by many people independently:

```
GROWTREE( $S$ )  
  if ( $y = 0$  for all  $\langle \mathbf{x}, y \rangle \in S$ ) return new leaf(0)  
  else if ( $y = 1$  for all  $\langle \mathbf{x}, y \rangle \in S$ ) return new leaf(1)  
  else  
    choose best attribute  $x_j$   
     $S_0 =$  all  $\langle \mathbf{x}, y \rangle \in S$  with  $x_j = 0$ ;  
     $S_1 =$  all  $\langle \mathbf{x}, y \rangle \in S$  with  $x_j = 1$ ;  
    return new node( $x_j$ , GROWTREE( $S_0$ ), GROWTREE( $S_1$ ))
```

## Choosing the Best Attribute

One way to choose the best attribute is to perform a 1-step lookahead search and choose the attribute that gives the lowest error rate on the training data.

**CHOOSEBESTATTRIBUTE(S)**

choose  $j$  to minimize  $J_j$ , computed as follows:

$S_0$  = all  $\langle \mathbf{x}, y \rangle \in S$  with  $x_j = 0$ ;

$S_1$  = all  $\langle \mathbf{x}, y \rangle \in S$  with  $x_j = 1$ ;

$y_0$  = the most common value of  $y$  in  $S_0$

$y_1$  = the most common value of  $y$  in  $S_1$

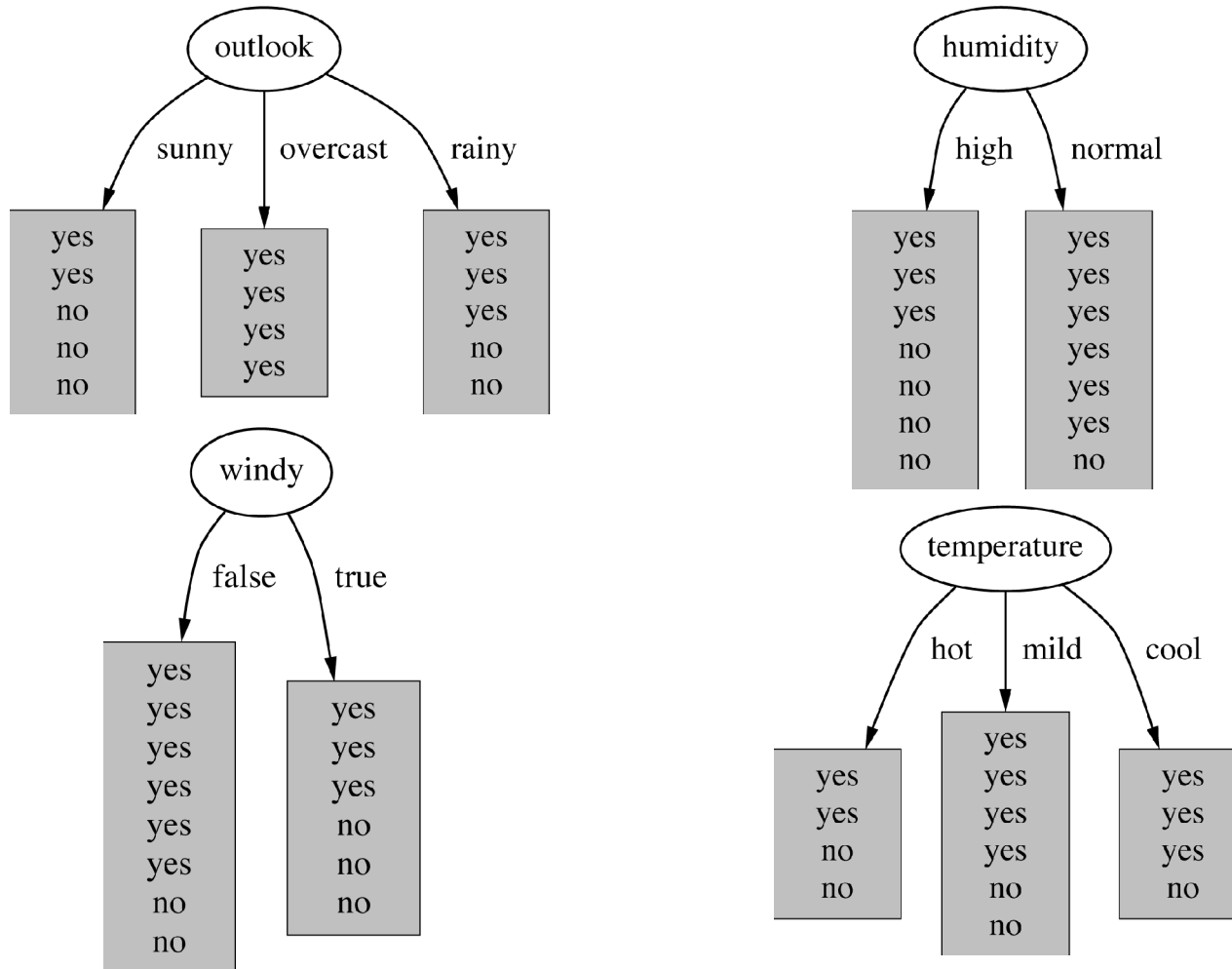
$J_0$  = number of examples  $\langle \mathbf{x}, y \rangle \in S_0$  with  $y \neq y_0$

$J_1$  = number of examples  $\langle \mathbf{x}, y \rangle \in S_1$  with  $y \neq y_1$

$J_j = J_0 + J_1$  (total errors if we split on this feature)

**return**  $j$

# Which attribute to select?

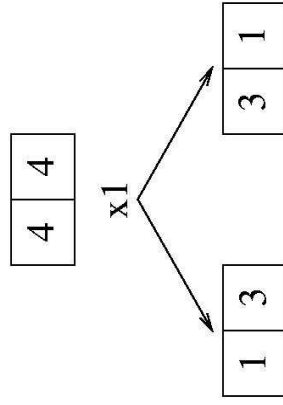


# A criterion for attribute selection

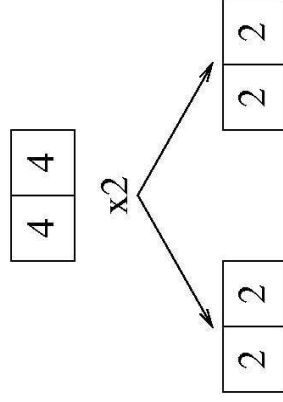
- Which is the best attribute?
  - The one which will result in the smallest tree
  - Heuristic: choose the attribute that produces the “purest” nodes
- Need a good measure of purity!
  - Maximal when?
  - Minimal when?

# Choosing the Best Attribute—An Example

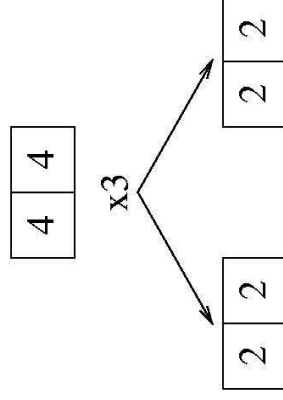
$x_1$	$x_2$	$x_3$	$y$
0	0	0	1
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	0



J=2



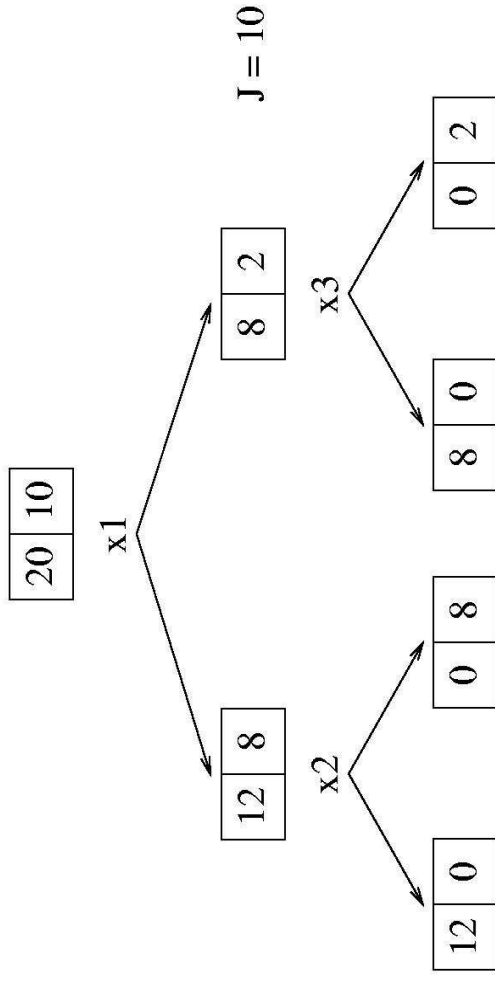
J=4



J=4

### Choosing the Best Attribute (3)

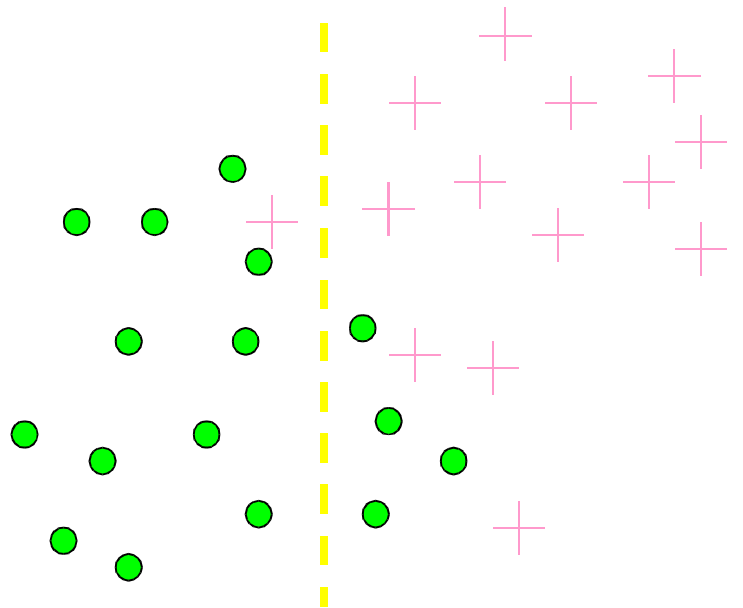
Unfortunately, this measure does not always work well, because it does not detect cases where we are making “progress” toward a good tree.



# Information Gain

Which test is more informative?

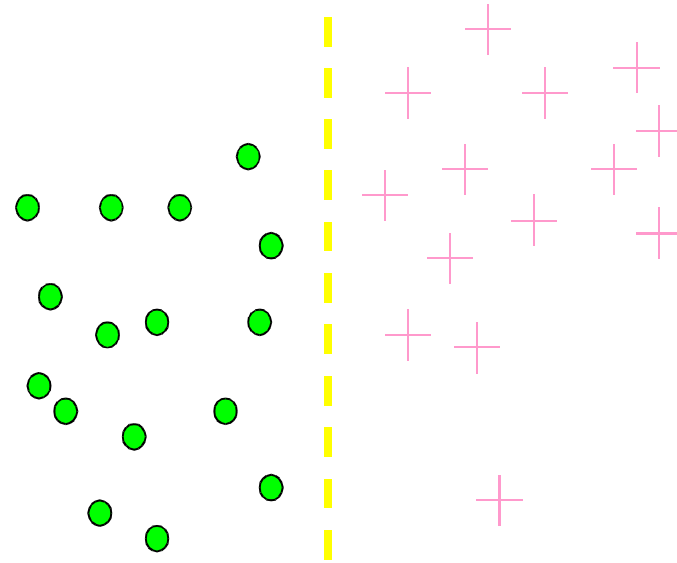
**Split over whether  
Balance exceeds 50K**



Less or equal 50K

Over 50K

**Split over whether  
applicant is employed**



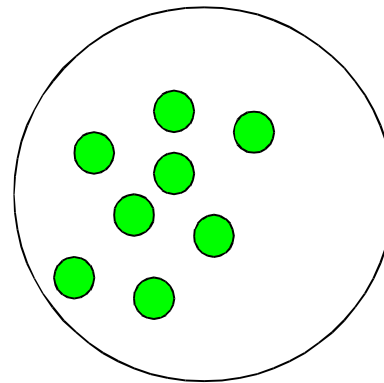
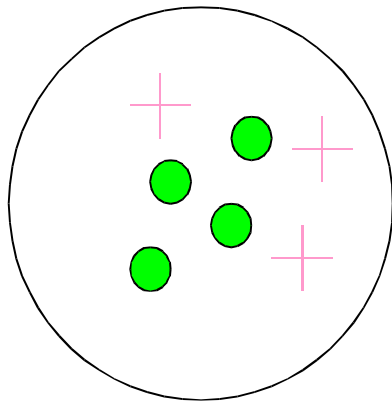
Unemployed

Employed

# Information Gain

## Impurity / Entropy (informal)

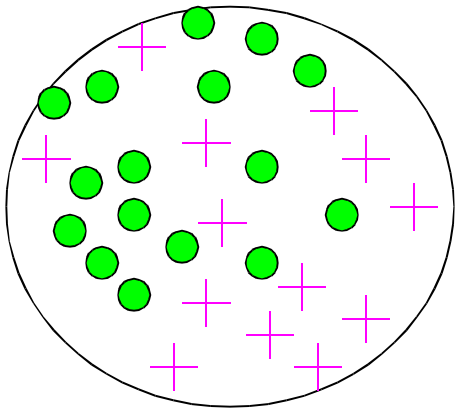
- Measures the level of **impurity** in a group of examples



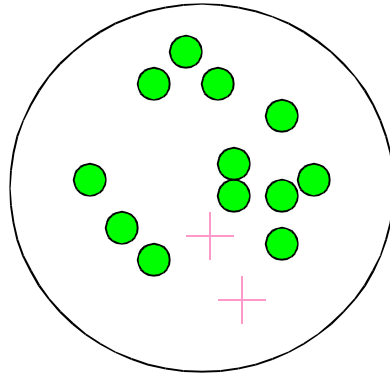


# Impurity

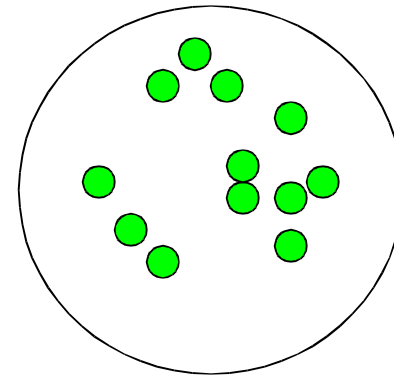
**Very impure group**



**Less impure**

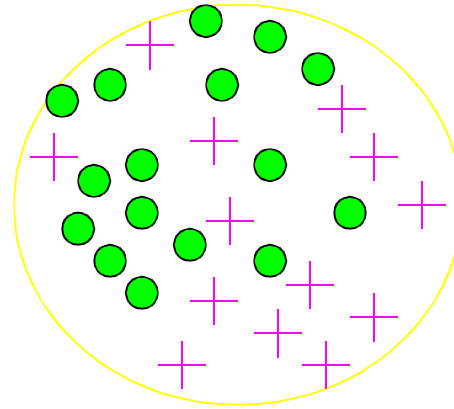


**Minimum  
impurity**



# Calculating Impurity

- Impurity =  $\sum -p_i \log_2 p_i$   
Pi is proportion of class i



**When examples can belong to one of two classes:**

**What is the worst case of impurity?**

# 2-class Cases:

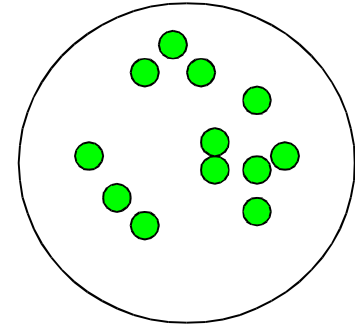
- What is the impurity of a group in which all examples belong to the same class?

– Impurity =  $-1 \log_2 1 = 0$

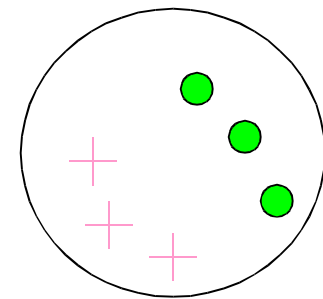
- What is the impurity of a group with 50% in either class?

– Impurity =  $-0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$

**Minimum  
impurity**



**Maximum  
impurity**

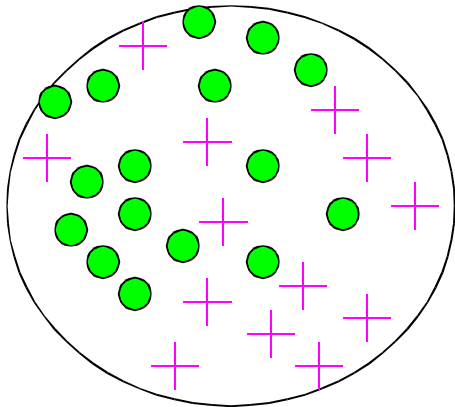


# Calculating Information Gain

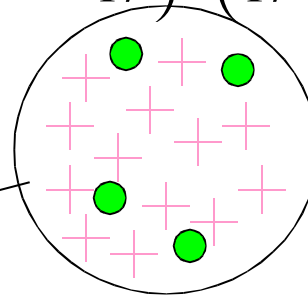
**Information Gain** = Impurity (parent) - [Impurity (children)]

$$\text{impurity} = -\left(\frac{13}{17} \cdot \log_2 \frac{13}{17}\right) - \left(\frac{4}{17} \cdot \log_2 \frac{4}{17}\right) = 0.787$$

Entire population (30 instances)

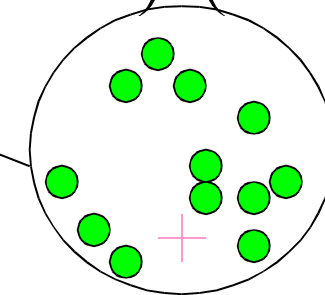


$$\text{impurity} = -\left(\frac{1}{13} \cdot \log_2 \frac{1}{13}\right) - \left(\frac{12}{13} \cdot \log_2 \frac{12}{13}\right) = 0.391$$



17 instances

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



13 instances

$$\text{(Weighted) Average Impurity of Children} = \left(\frac{17}{30} \cdot 0.787\right) + \left(\frac{13}{30} \cdot 0.391\right) = 0.615$$

$$\text{Information Gain} = 0.996 - 0.615 = 0.38$$

## Non-Boolean Features

- **Features with multiple discrete values**

Construct a multiway split?

Test for one value versus all of the others?

Group the values into two disjoint subsets?

- **Real-valued features**

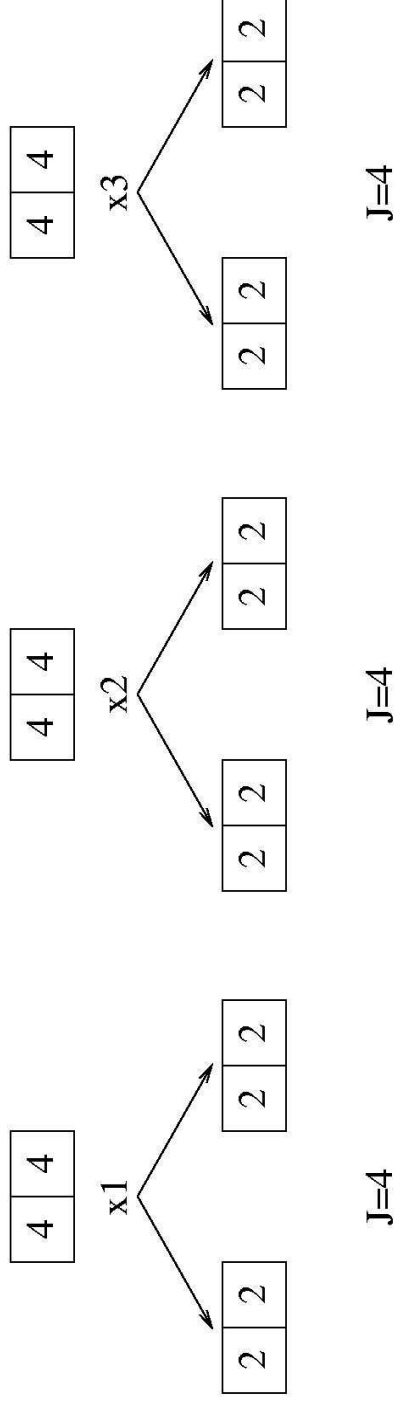
Consider a threshold split using each observed value of the feature.

Whichever method is used, the mutual information can be computed to choose the best split.

## Learning Parity with Noise

When learning exclusive-or (2-bit parity), all splits look equally good. If extra random boolean features are included, they also look equally good. Hence, decision tree algorithms cannot distinguish random noisy features from parity features.

$x_1$	$x_2$	$x_3$	$y$
0	0	0	0
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	1
1	0	1	1
1	1	0	0
1	1	1	0



## Unknown Attribute Values

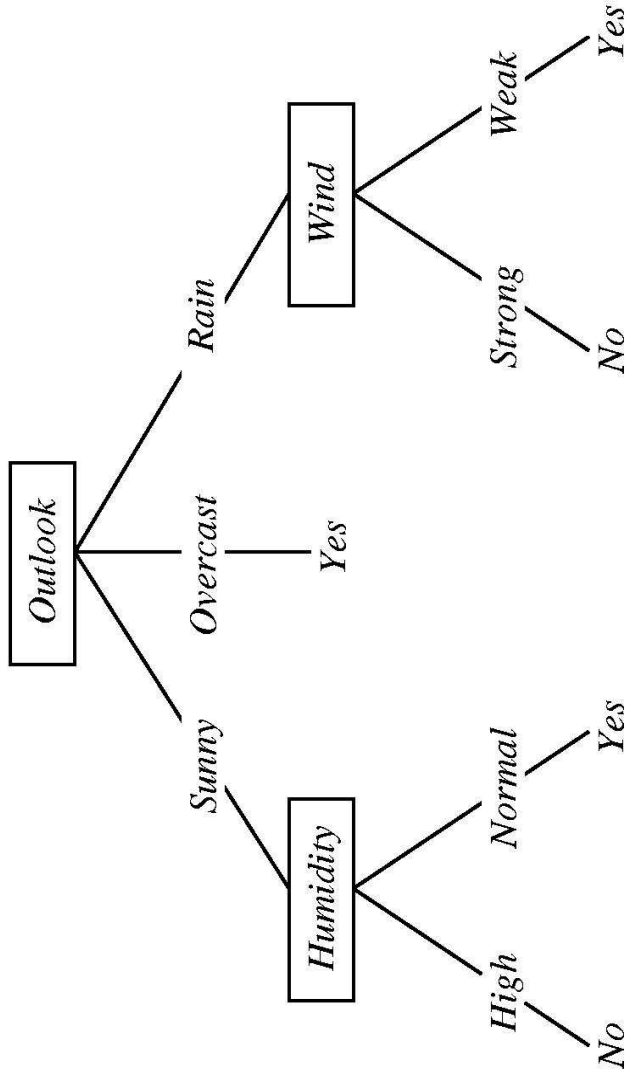
What if some examples are missing values of  $A$ ?

Use training example anyway, sort through tree

- If node  $n$  tests  $A$ , assign most common value of  $A$  among other examples sorted to node  $n$
- Assign most common value of  $A$  among other examples with same target value
- Assign probability  $p_i$  to each possible value  $v_i$  of  $A$   
Assign fraction  $p_i$  of example to each descendant in tree

Classify new examples in same fashion

# Overfitting in Decision Trees



Consider adding a noisy training example:

*Sunny, Hot, Normal, Strong, PlayTennis=No*

What effect on tree?



# Overfitting

Consider error of hypothesis  $h$  over

- training data:  $error_{train}(h)$
- entire distribution  $\mathcal{D}$  of data:  $error_{\mathcal{D}}(h)$

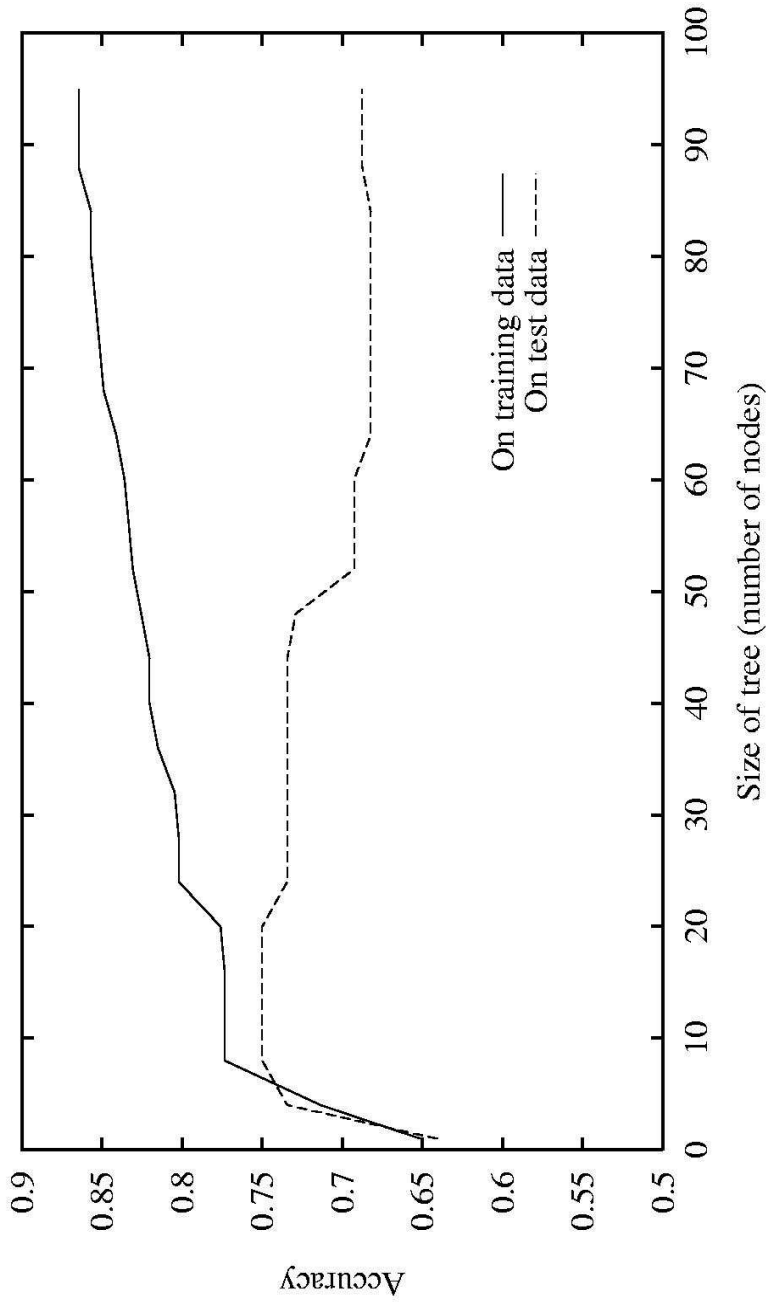
Hypothesis  $h \in H$  **overfits** training data if there is an alternative hypothesis  $h' \in H$  such that

$$error_{train}(h) < error_{train}(h')$$

and

$$error_{\mathcal{D}}(h) > error_{\mathcal{D}}(h')$$

# Overfitting in Decision Tree Learning



# Avoiding Overfitting

How can we avoid overfitting?

- Stop growing when data split not statistically significant
- Grow full tree, then post-prune

How to select “best” tree:

- Measure performance over training data
- Measure performance over separate validation data set
- Add complexity penalty to performance measure

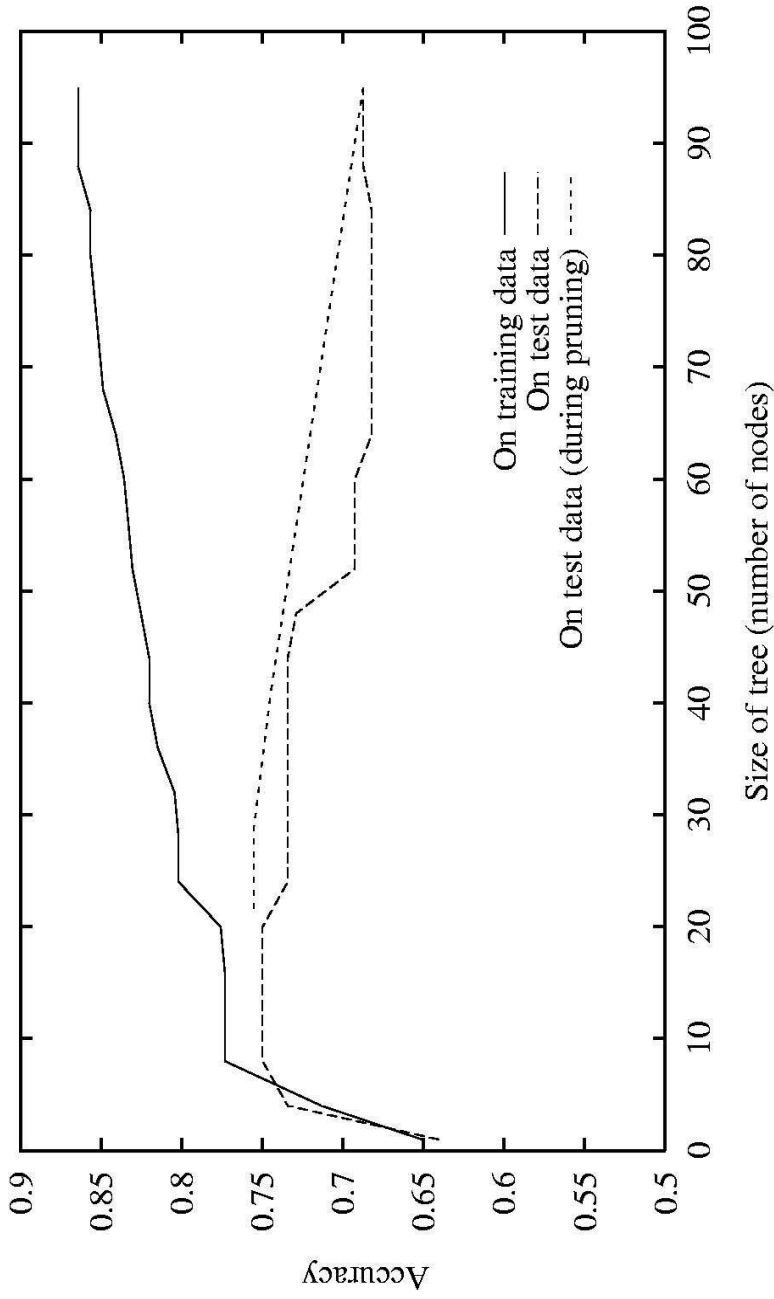
# Reduced-Error Pruning

Split data into *training* and *validation* set

Do until further pruning is harmful:

1. Evaluate impact on *validation* set of pruning each possible node (plus those below it)
2. Greedily remove the one that most improves *validation* set accuracy

# Effect of Reduced-Error Pruning

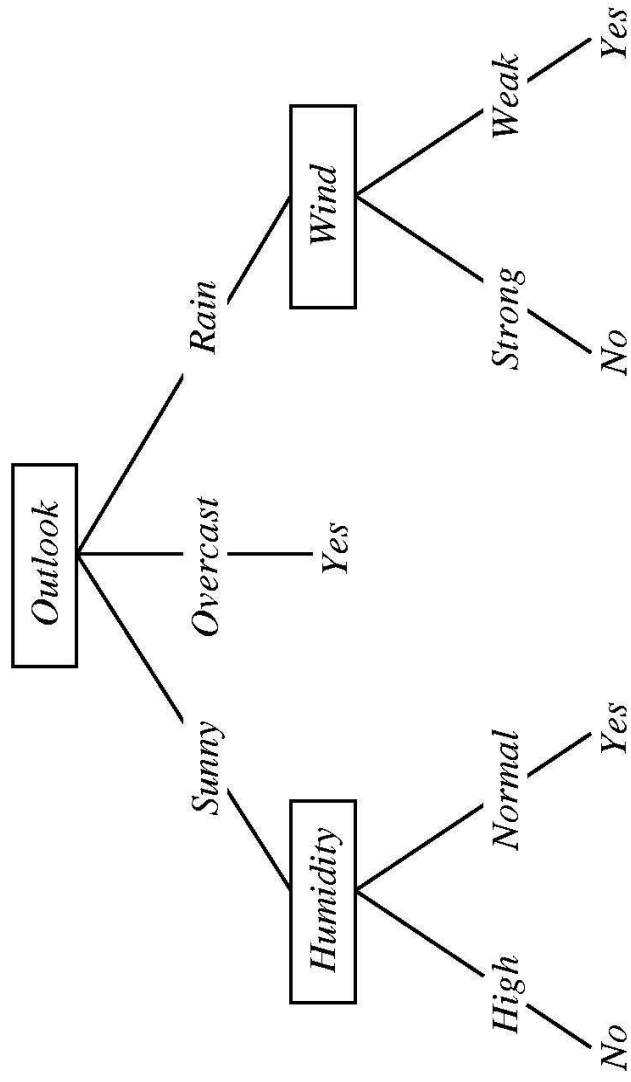


## Rule Post-Pruning

1. Convert tree to equivalent set of rules
2. Prune each rule independently of others
3. Sort final rules into desired sequence for use

Perhaps most frequently used method (e.g., C4.5)

# Converting A Tree to Rules



IF     (*Outlook = Sunny*) AND (*Humidity = High*)  
THEN  *PlayTennis = No*

IF     (*Outlook = Sunny*) AND (*Humidity = Normal*)  
THEN  *PlayTennis = Yes*

...



## Scaling Up

- ID3, C4.5, etc. assume data fits in main memory  
(OK for up to hundreds of thousands of examples)
- SPRINT, SLIQ: multiple sequential scans of data  
(OK for up to millions of examples)
- VFDT: at most one sequential scan  
(OK for up to billions of examples)

# Decision Trees: Summary

- Representation = decision trees
- Bias = preference for small decision trees
- Search algorithm =
- Heuristic function = information gain
- Overfitting and pruning