

# CLASSIFICATION DECISION TREES



# Classification

- ❖ Given a collection of records (*training set*)
  - ✓ Each record contains a set of *attributes*, one of the attributes is the *class*.
- ❖ Find a *model* for class attribute as a function of the values of other attributes.
- ❖ Goal: previously unseen records should be assigned a class as accurately as possible.
  - ✓ A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

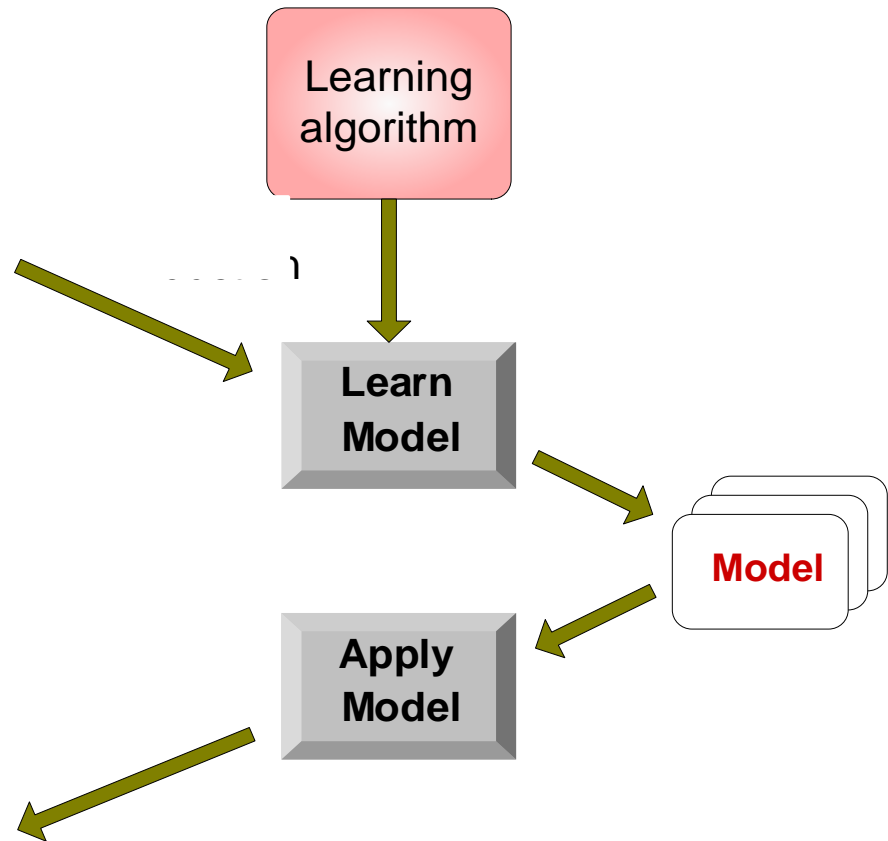
# Classification

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



# Classification

		Predicted Class	
		<i>Class</i> = 1	<i>Class</i> = 0
Actual Class	<i>Class</i> = 1	$f_{11}$	$f_{10}$
	<i>Class</i> = 0	$f_{01}$	$f_{00}$

Confusion Matrix

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

$$\text{Error rate} = \frac{\text{Number of wrong predictions}}{\text{Total number of predictions}} = \frac{f_{10} + f_{01}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

# Decision Trees

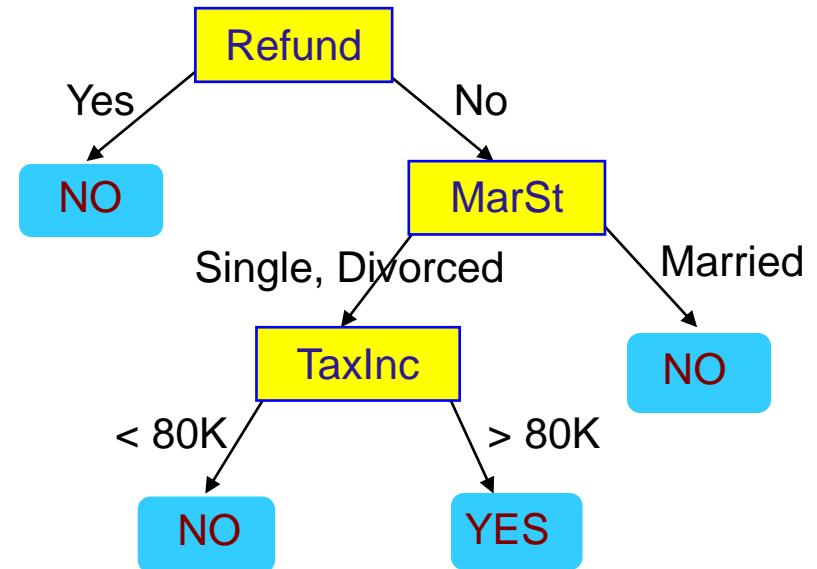
- ✓ solve a classification problem by asking a series of carefully crafted questions about the attributes of the test record
- ✓ series of questions and their possible answers can be organized in the form of a decision Tree
- ✓ nodes and directed edges
  - ❖ **root node**
  - ❖ **Internal nodes**
  - ❖ **Leaf or terminal nodes**

# Decision Trees

categorical  
categorical  
continuous  
class

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Training Data



Model: Decision Tree

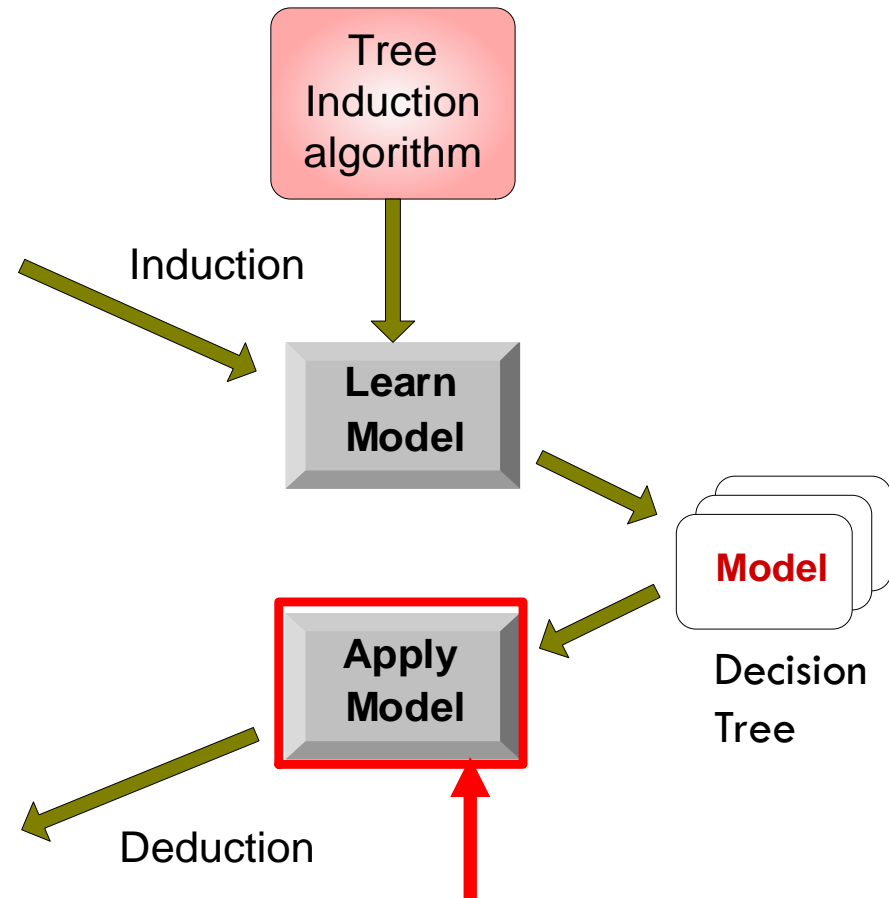
# Decision Tree Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set

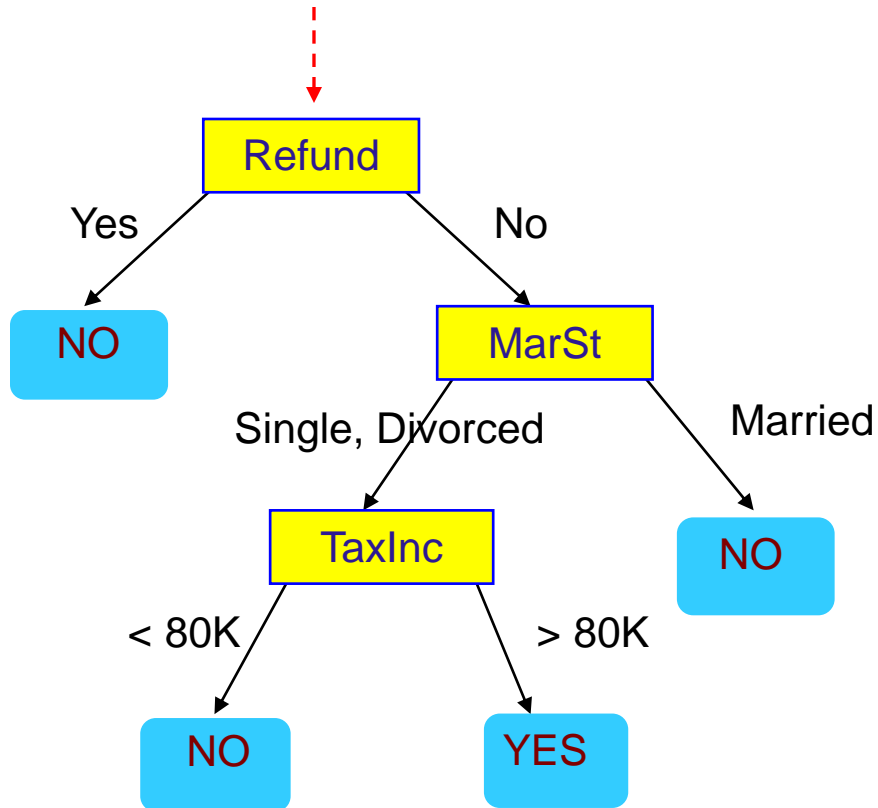


# Apply Model to Test Data

## Test Data

Start from the root of tree.

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

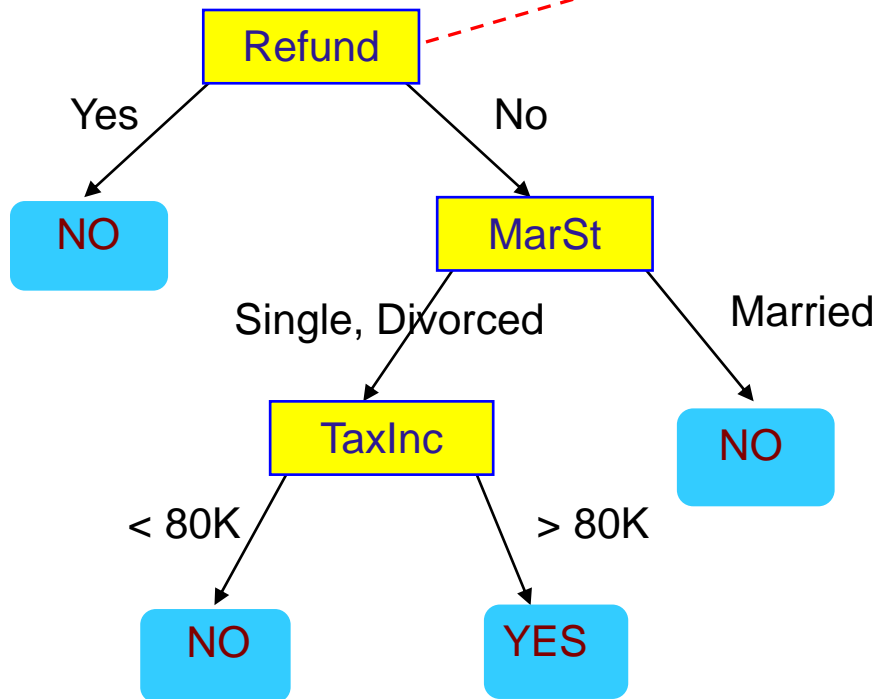




# Apply Model to Test Data

## Test Data

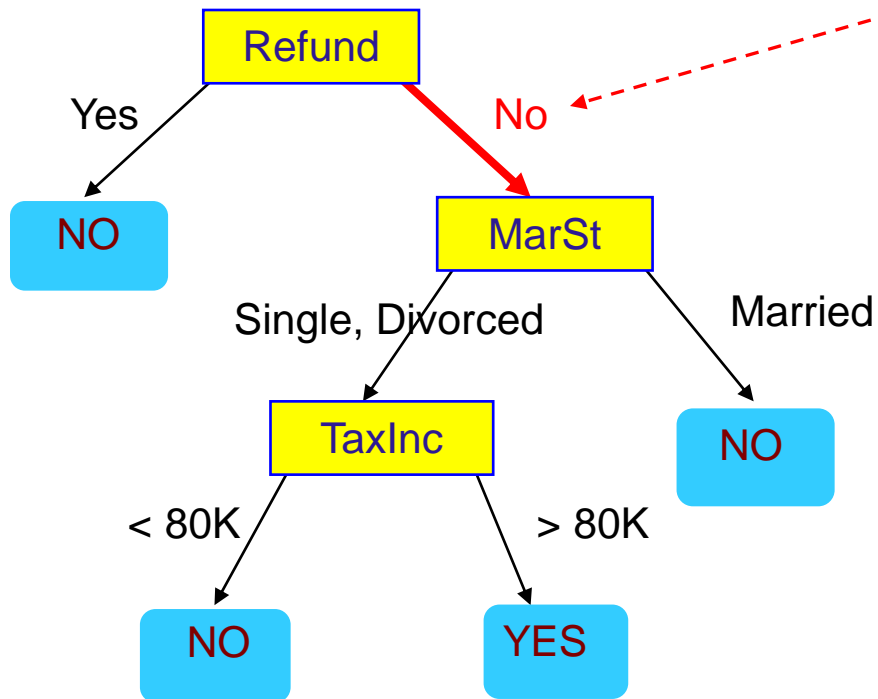
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



# Apply Model to Test Data

## Test Data

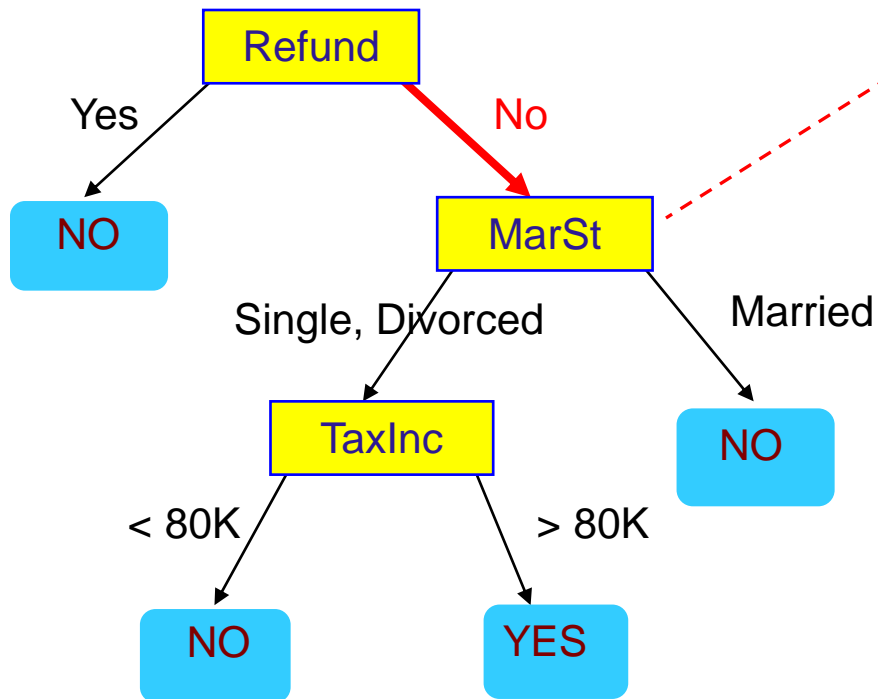
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



# Apply Model to Test Data

## Test Data

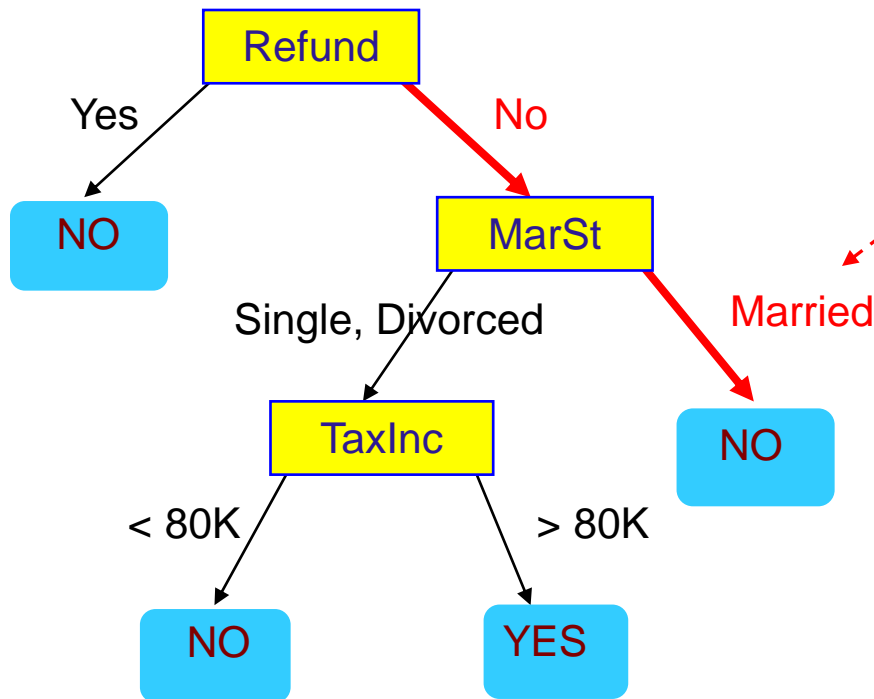
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



# Apply Model to Test Data

## Test Data

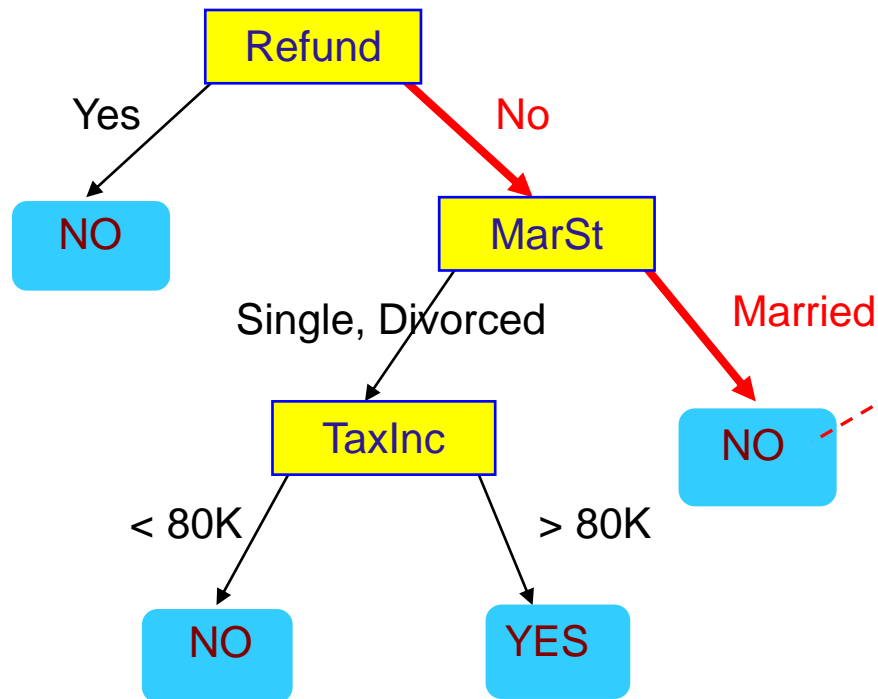
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



# Apply Model to Test Data

## Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



Assign Cheat to "No"

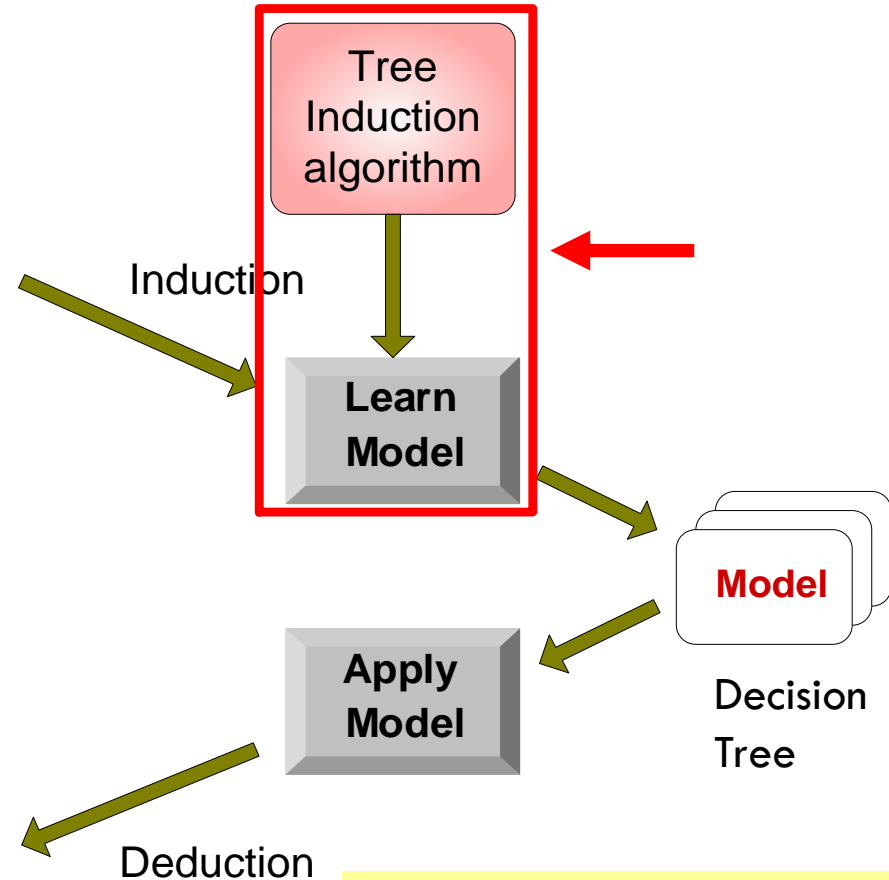
# Decision Tree Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



How to Build a Decision Tree

# Decision Trees

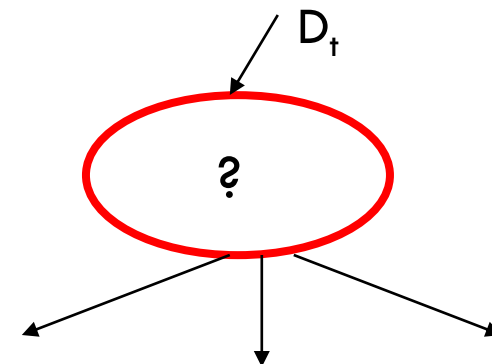
- ❖ many decision trees
- ❖ finding the optimal tree is computationally infeasible
- ❖ efficient algorithms to induce a reasonably accurate, albeit suboptimal, decision tree in a reasonable amount of time
- ❖ Employ a **greedy strategy**
- ❖ grows a decision tree by making a series of locally optimum decisions about which attribute to use for partitioning the data

## Hunt's Algorithm

# General Structure of Hunt's Algorithm

- Let  $D_t$  be the set of training records that reach a node  $t$
- General Procedure:
  - ▣ If  $D_t$  contains records that belong the same class  $y_t$ , then  $t$  is a leaf node labeled as  $y_t$
  - ▣ If  $D_t$  is an empty set, then  $t$  is a leaf node labeled by the default class,  $y_d$
  - ▣ If  $D_t$  contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

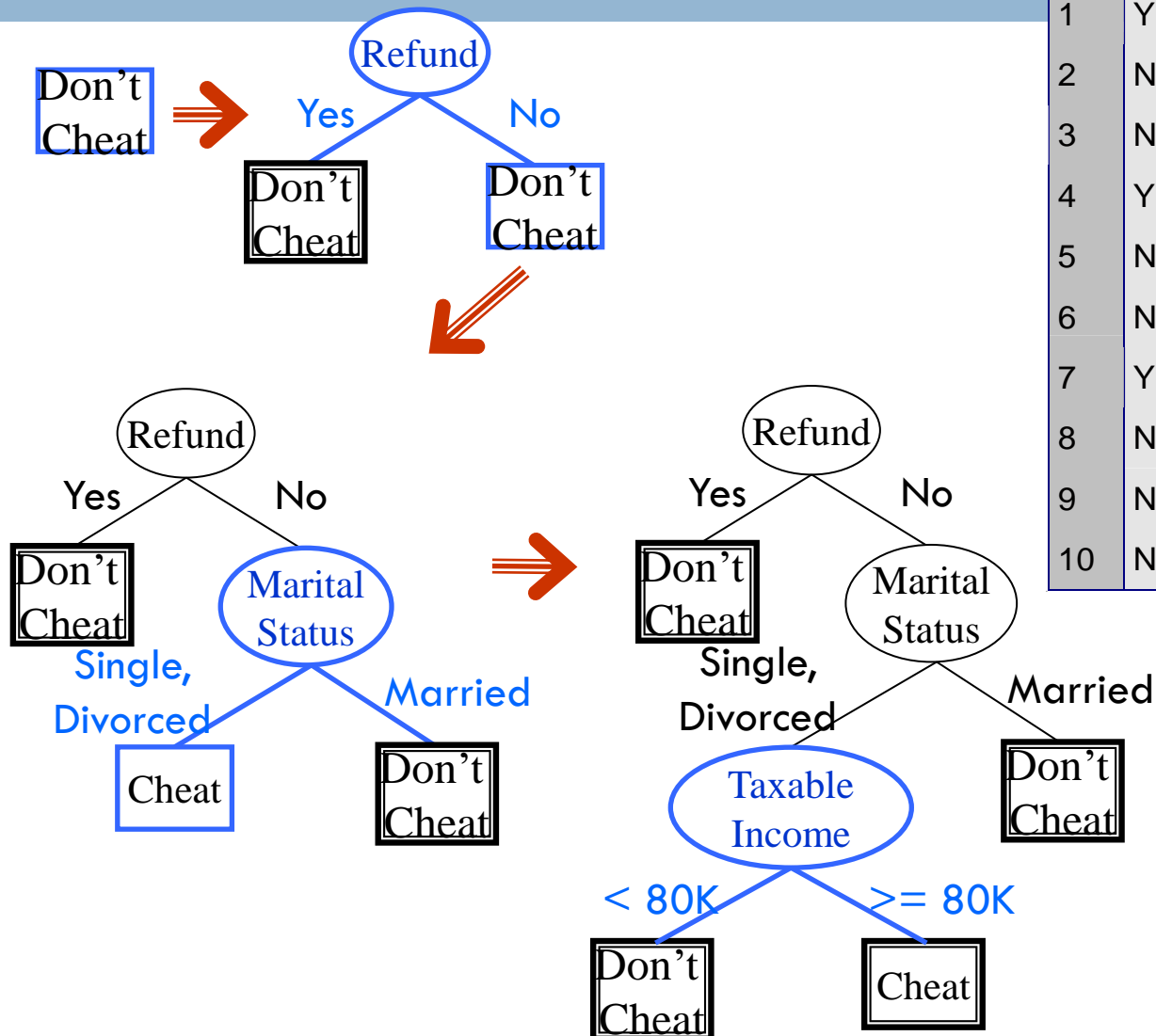
<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes





# Hunt's Algorithm

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



# Tree Induction

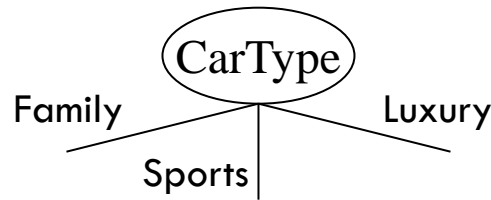
- Greedy strategy.
  - ▣ Split the records based on an attribute test that optimizes certain criterion.
  
- Issues
  - ▣ Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - ▣ Determine when to stop splitting

# How to Specify Test Condition?

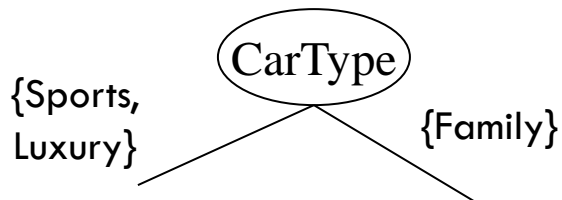
- Depends on attribute types
  - ▣ Nominal
  - ▣ Ordinal
  - ▣ Continuous
  
- Depends on number of ways to split
  - ▣ 2-way split
  - ▣ Multi-way split

# Splitting Based on Nominal Attributes

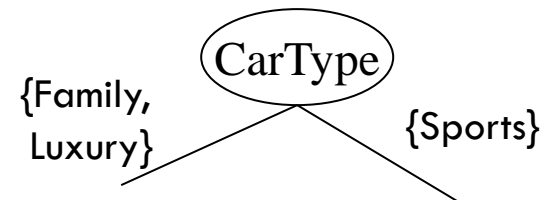
- **Multi-way split:** Use as many partitions as distinct values.



- **Binary split:** Divides values into two subsets.  
Need to find optimal partitioning.

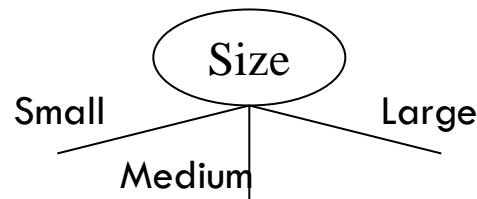


OR

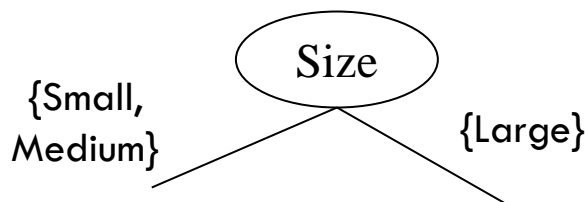


# Splitting Based on Ordinal Attributes

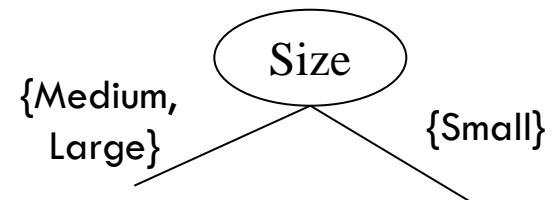
- **Multi-way split:** Use as many partitions as distinct values.



- **Binary split:** Divides values into two subsets.  
Need to find optimal partitioning.



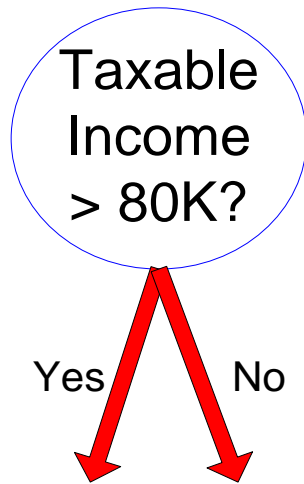
OR



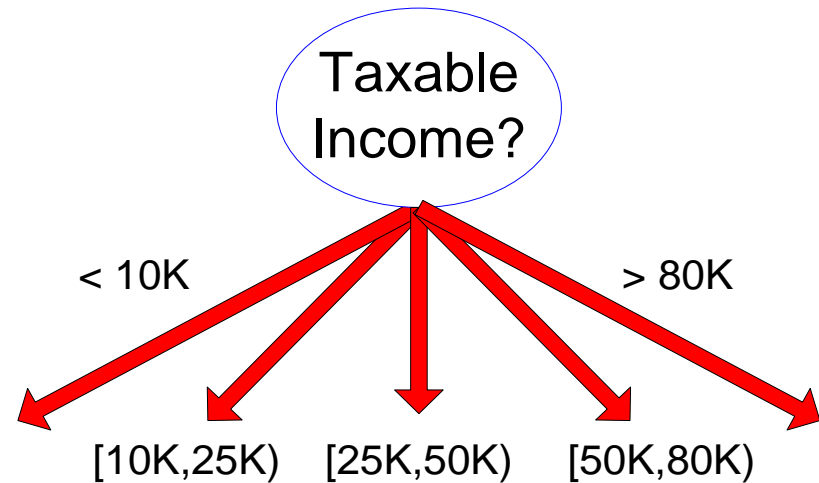
# Splitting Based on Continuous Attributes

- Different ways of handling
  - ▣ **Discretization** to form an ordinal categorical attribute
  - ▣ **Binary Decision**:  $(A < v)$  or  $(A \geq v)$ 
    - consider all possible splits and finds the best cut

# Splitting Based on Continuous Attributes



(i) Binary split



(ii) Multi-way split

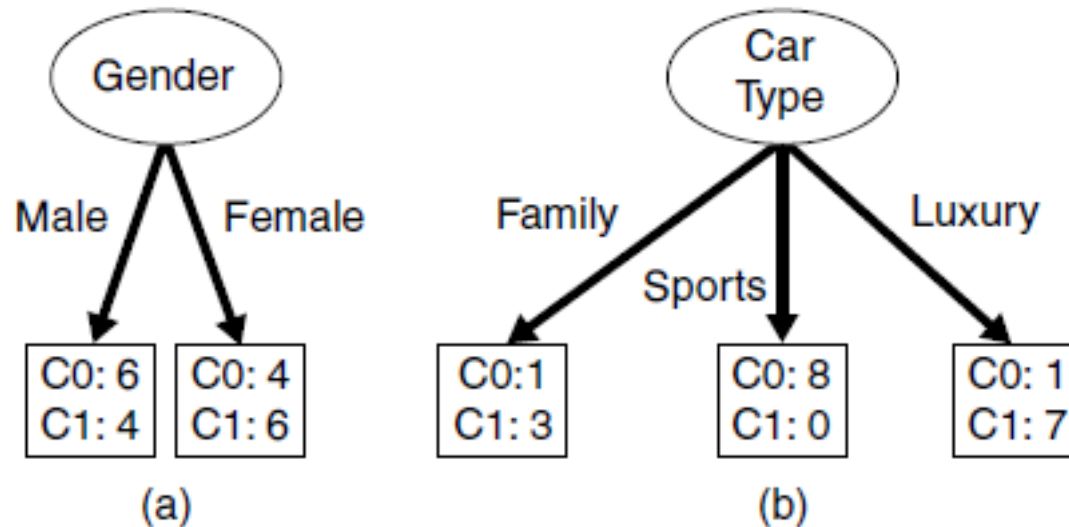
# Tree Induction

- Greedy strategy.
  - ▣ Split the records based on an attribute test that optimizes certain criterion.
  
- Issues
  - ▣ Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - ▣ Determine when to stop splitting



# How to determine the Best Split

Before Splitting: 10 records of class 0,  
10 records of class 1



Which test condition is the best?

# How to determine the Best Split

- Greedy approach:
  - ▣ Nodes with **homogeneous** class distribution are preferred
- Need a measure of node impurity:

C0: 5  
C1: 5

Non-homogeneous,  
High degree of impurity

C0: 9  
C1: 1

Homogeneous,  
Low degree of impurity

- Gini Index
- Entropy
- Misclassification on error

# Measure of Impurity: GINI

- Gini Index for a given node  $t$  :

$$GINI(t) = 1 - \sum_j [p(j|t)]^2$$

(NOTE:  $p(j|t)$  is the relative frequency of class  $j$  at node  $t$ ).

- ▣ Maximum  $(1 - 1/n_c)$  when records are equally distributed among all classes, implying least interesting information
- ▣ Minimum (0) when all records belong to one class, implying most interesting information

# Examples for computing GINI

$$GINI(t) = 1 - \sum_j [p(j | t)]^2$$

C1	<b>0</b>
C2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Gini = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

C1	<b>1</b>
C2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Gini = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

C1	<b>2</b>
C2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$Gini = 1 - (2/6)^2 - (4/6)^2 = 0.444$$

# Splitting Based on GINI

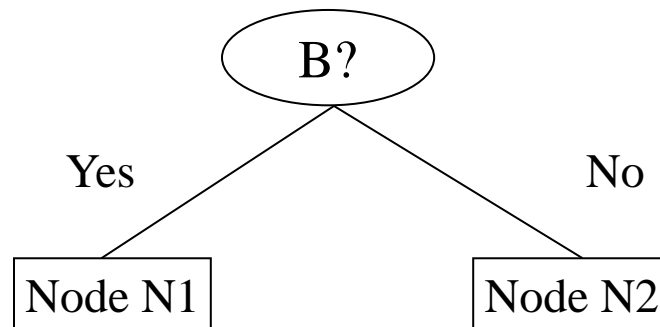
- Used in CART, SLIQ, SPRINT.
- When a node  $p$  is split into  $k$  partitions (children), the quality of split is computed as,

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

where,  $n_i$  = number of records at child  $i$ ,  
 $n$  = number of records at node  $p$ .

# Binary Attributes: Computing GINI Index

- Splits into two partitions

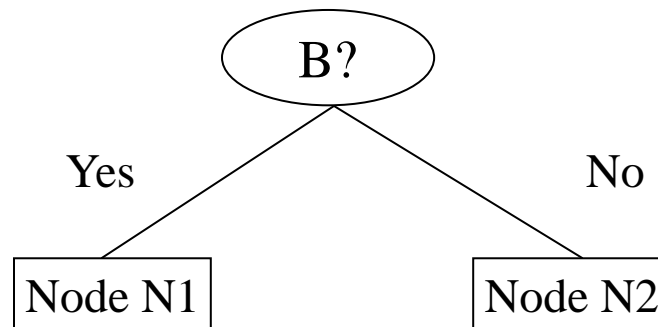


	N1	N2
C1	5	1
C2	2	4
Gini=0.333		

	Parent
C1	6
C2	6
Gini = 0.500	

# Binary Attributes: Computing GINI Index

- Splits into two partitions



$$\begin{aligned} \text{Gini}(N1) &= 1 - (5/7)^2 - (2/7)^2 \\ &= 0.194 \end{aligned}$$

$$\begin{aligned} \text{Gini}(N2) &= 1 - (1/5)^2 - (4/5)^2 \\ &= 0.528 \end{aligned}$$

	<b>N1</b>	<b>N2</b>
C1	<b>5</b>	<b>1</b>
C2	<b>2</b>	<b>4</b>
<b>Gini=0.333</b>		

	<b>Parent</b>
C1	<b>6</b>
C2	<b>6</b>
<b>Gini = 0.500</b>	

$$\begin{aligned} \text{Gini(Children)} &= 7/12 * 0.194 + \\ &\quad 5/12 * 0.528 \\ &= 0.333 \end{aligned}$$

# Categorical Attributes: Computing Gini Index

- For each distinct value, gather counts for each class in the dataset
- Use the count matrix to make decisions

Multi-way split

	CarType		
	Family	Sports	Luxury
C1	1	2	1
C2	4	1	1
Gini	0.393		

Two-way split  
(find best partition of values)

	CarType	
	{Sports, Luxury}	{Family}
C1	3	1
C2	2	4
Gini	0.400	

	CarType	
	{Sports}	{Family, Luxury}
C1	2	2
C2	1	5
Gini	0.419	



# Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
  - ▣ Sort the attribute on values
  - ▣ scan these values, each time updating the count matrix and computing gini index
  - ▣ Choose the split position that has the least gini index

Cheat	No		No		No		Yes		Yes		Yes		No		No		No		No			
→  →	Taxable Income																					
	60		70		75		85		90		95		100		120		125		220			
	55		65		72		80		87		92		97		110		122		172		230	
	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>
Yes	0	3	0	3	0	3	0	3	1	2	2	1	3	0	3	0	3	0	3	0	3	0
No	0	7	1	6	2	5	3	4	3	4	3	4	3	4	4	3	5	2	6	1	7	0
Gini	0.420		0.400		0.375		0.343		0.417		0.400		<u>0.300</u>		0.343		0.375		0.400		0.420	

## Alternative Splitting Criteria based on INFO

- Entropy at a given node  $t$ :

$$Entropy(t) = -\sum_j p(j | t) \log p(j | t)$$

(NOTE:  $p(j | t)$  is the relative frequency of class  $j$  at node  $t$ ).

- ▣ Measures homogeneity of a node.
  - Maximum ( $\log n_c$ ) when records are equally distributed among all classes implying least information
  - Minimum (0.0) when all records belong to one class, implying most information
- ▣ Entropy based computations are similar to the GINI index computations

# Examples for computing Entropy

$$Entropy(t) = -\sum_j p(j|t) \log_2 p(j|t)$$

C1	<b>0</b>
C2	<b>6</b>

C1	<b>1</b>
C2	<b>5</b>

C1	<b>2</b>
C2	<b>4</b>

# Examples for computing Entropy

$$Entropy(t) = -\sum_j p(j|t) \log_2 p(j|t)$$

C1	<b>0</b>
C2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Entropy = -0 \log 0 - 1 \log 1 = -0 - 0 = 0$$

C1	<b>1</b>
C2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Entropy = -(1/6) \log_2 (1/6) - (5/6) \log_2 (5/6) = 0.65$$

C1	<b>2</b>
C2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

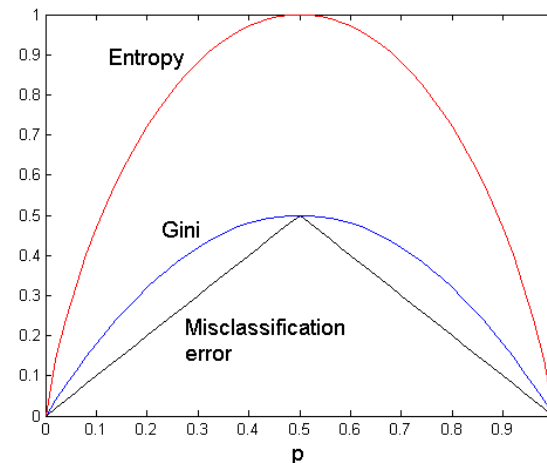
$$Entropy = -(2/6) \log_2 (2/6) - (4/6) \log_2 (4/6) = 0.92$$

# Splitting Criteria based on Classification Error

- Classification error at a node  $t$  :

$$Error(t) = 1 - \max_i P(i | t)$$

- Measures misclassification error made by a node.
  - Maximum  $(1 - 1/n_c)$  when records are equally distributed among all classes, implying least interesting information
  - Minimum (0.0) when all records belong to one class, implying most interesting information



# Gain

gain,  $\Delta$ , is a criterion that can be used to determine the goodness of a split

Entropy:  
Information  
Gain

$$\Delta = I(\text{parent}) - \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j)$$

$I(\cdot)$  is the impurity measure

$N$  is the total number of records at the parent node

$k$  is the number of attribute values

$N(v_j)$  is the number of records associated with the child node,  $v_j$ .

Decision tree induction algorithms often choose a test condition that maximizes the gain

# Gain ratio

$$\textit{GainRatio}_{split} = \frac{\textit{GAIN}_{Split}}{\textit{SplitINFO}}$$

$$\textit{SplitINFO} = -\sum_{i=1}^k \frac{n_i}{n} \log \frac{n_i}{n}$$

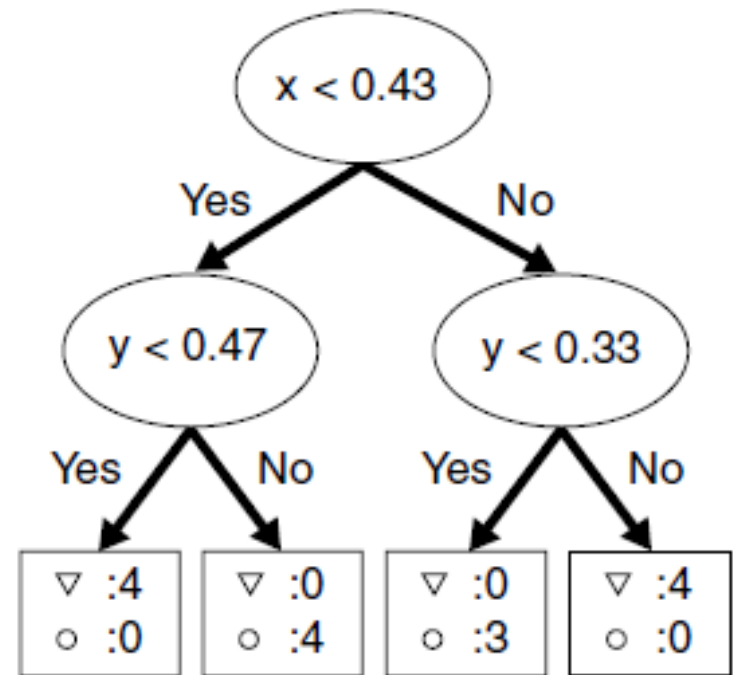
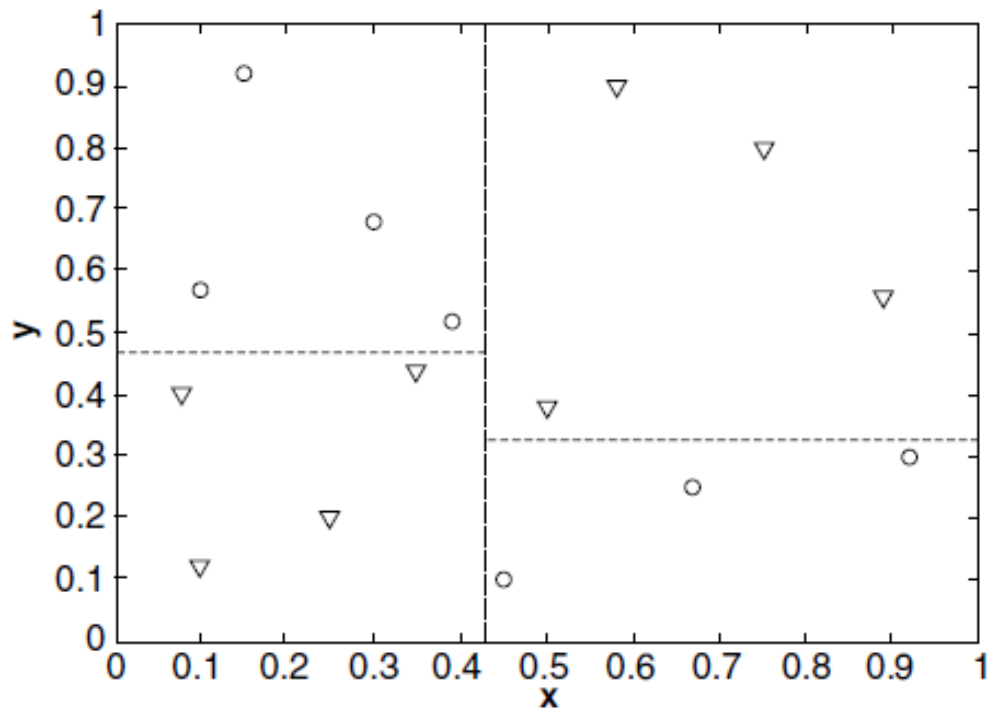
Parent Node, p is split into k partitions  
 $n_i$  is the number of records in partition i  
entropy of the partitioning

# Characteristics

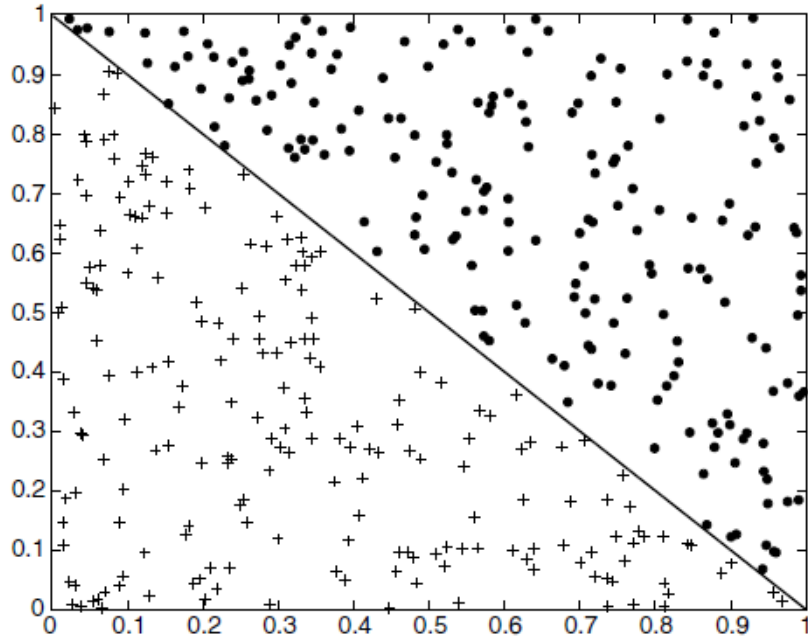
- ✓ Decision tree induction is a nonparametric approach for building classification Models
- ✓ Finding an optimal decision tree is a hard problem
- ✓ greedy approach
- ✓ relatively easy to interpret
- ✓ choice of impurity measure has little effect on the performance of decision tree induction algorithms
- ✓ using only a single attribute at a time
- ✓ partitioning the attribute space into disjoint regions until each region contains records of the same class
- ✓ parallel to the “coordinate axes.”



# Characteristics



# Characteristics



oblique decision tree

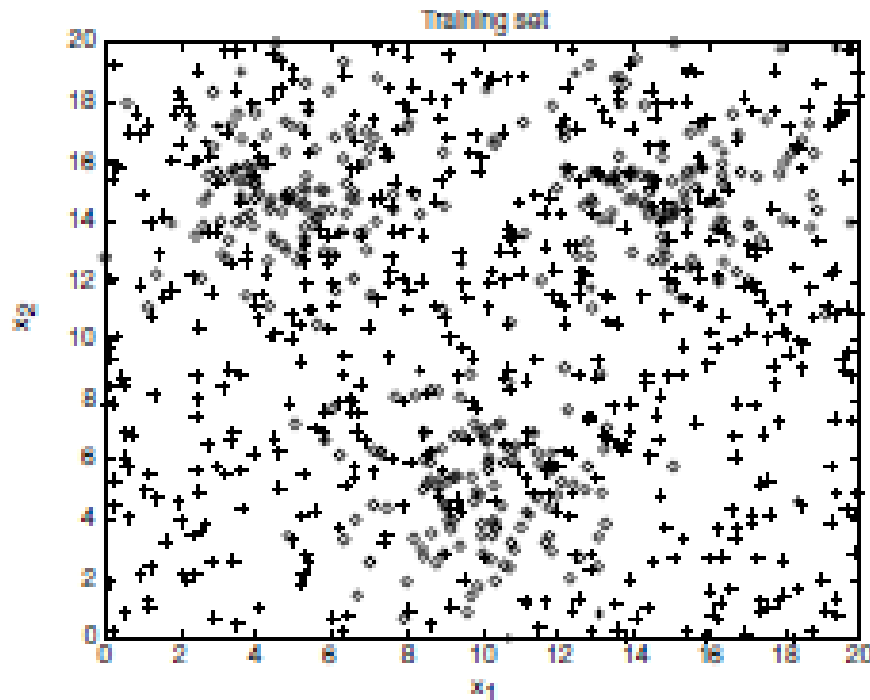
$$x + y < 1$$

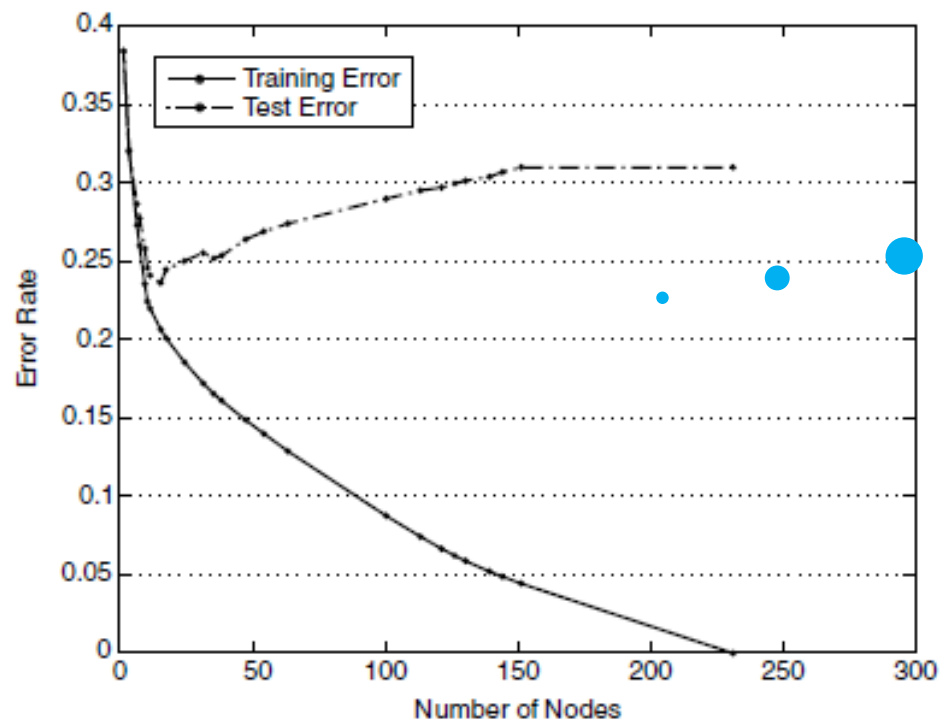
# Model Overfitting

Training errors

Test errors (generalization)

good classification model must not only fit training data well, it must also accurately classify records it has never

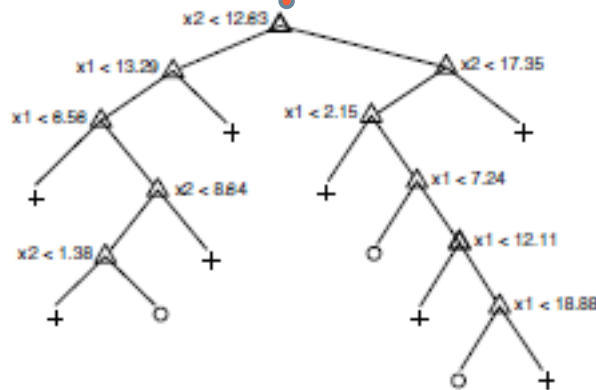




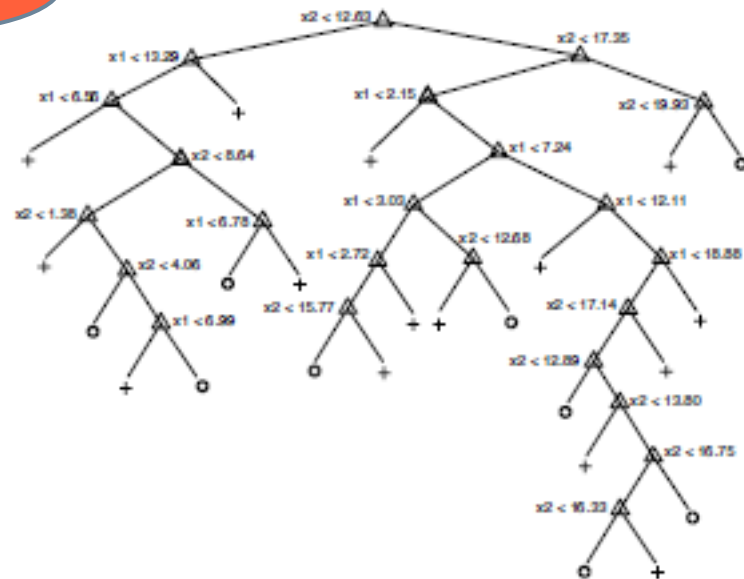
Underfitting  
and Overfitting

# Model Overfitting

a higher training error rate,  
but a lower test error rate



(a) Decision tree with 11 leaf nodes.



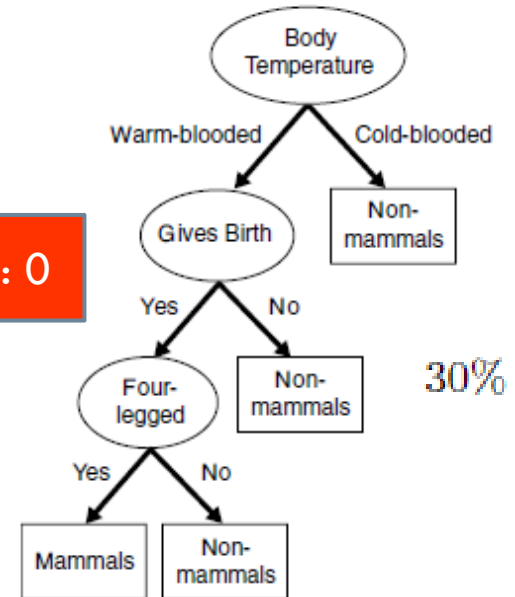
(b) Decision tree with 24 leaf nodes.

# Overfitting Due to Presence of Noise

## Train

Name	Body Temperature	Gives Birth	Four-legged	Hibernates	Class Label
porcupine	warm-blooded	yes	yes	yes	yes
cat	warm-blooded	yes	yes	no	yes
bat	warm-blooded	yes	no	yes	no*
whale	warm-blooded	yes	no	no	no*
salamander	cold-blooded	no	yes	yes	no
komodo dragon	cold-blooded	no	yes	no	no
python	cold-blooded	no	no	yes	no
salmon	cold-blooded	no	no	no	no
eagle	warm-blooded	no	no	no	no
guppy	cold-blooded	yes	no	no	no

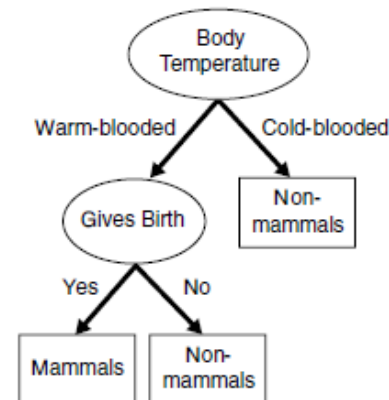
Train error: 0



(a) Model M1

## Test

Name	Body Temperature	Gives Birth	Four-legged	Hibernates	Class Label
human	warm-blooded	yes	no	no	yes
pigeon	warm-blooded	no	no	no	no
elephant	warm-blooded	yes	yes	no	yes
leopard shark	cold-blooded	yes	no	no	no
turtle	cold-blooded	no	yes	no	no
penguin	cold-blooded	no	no	no	no
eel	cold-blooded	no	no	no	no
dolphin	warm-blooded	yes	no	no	yes
spiny anteater	warm-blooded	no	yes	yes	yes
gila monster	cold-blooded	no	yes	yes	no

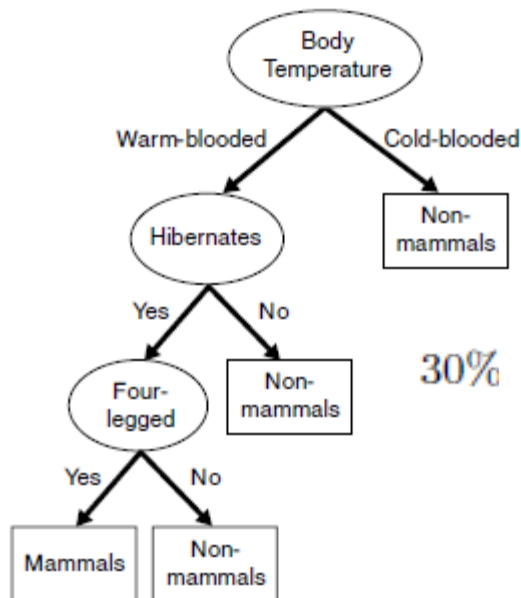


Train error:  
20%

error rate (10%)

# small number of training records

Name	Body Temperature	Gives Birth	Four-legged	Hibernates	Class Label
salamander	cold-blooded	no	yes	yes	no
guppy	cold-blooded	yes	no	no	no
eagle	warm-blooded	no	no	no	no
poorwill	warm-blooded	no	no	yes	no
platypus	warm-blooded	no	yes	yes	yes



# Decision Trees

## Prepruning :Early Stopping Rule

a more restrictive stopping condition

stop expanding a leaf node when the observed gain in impurity measure is low

## Post-pruning

decision tree is initially grown to its maximum size

tree-pruning step

replacing a subtree with a new leaf node



# Evaluating the Performance of a Classifier

accuracy or error rate computed from the test set can be used to compare different classifiers

class labels of test records must be known

## Holdout Method

1. labeled examples partitioned into two disjoint sets: training and the test sets
2. classification model is then induced from the training set
3. its performance is evaluated on the test set

- ✓ smaller training set size, larger variance of the model
- ✓ training set is too large, then the estimated accuracy computed from the smaller test set is less reliable

# Evaluating the Performance of a Classifier

## Random Subsampling

Repeated holdout

## Bootstrap

Sampling with replacement

## Cross-Validation

each record is used the same number of times for training and exactly once for testing

K-fold Cross-Validation