

Notes

- To preview next lecture:

- ▣ Check the lecture notes, if slides are not available:

<http://web.cse.ohio-state.edu/~sun.397/courses/au2017/cse5243-new.html>

- ▣ Check UIUC course on the same topic. All their slides are available:

<https://wiki.illinois.edu/wiki/display/cs412/2.+Course+Syllabus+and+Schedule>

- Extra readings beyond the lecture slides are important for related job hunting or research:

- ▣ Lecture notes (or, chapters from the text book Han et al.)

- ▣ Text book: <http://www.dataminingbook.info/pmwiki.php/Main/BookDownload>

CSE 5243 INTRO. TO DATA MINING

Classification (Basic Concepts)

Huan Sun, CSE@The Ohio State University

Classification: Basic Concepts

- Classification: Basic Concepts 
- Decision Tree Induction
- Bayes Classification Methods
- Model Evaluation and Selection
- Techniques to Improve Classification Accuracy: Ensemble Methods
- Summary

Supervised vs. Unsupervised Learning

- Supervised learning (classification)
 - ▣ Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations
 - ▣ New data is classified based on the training set

Supervised vs. Unsupervised Learning

- **Supervised learning (classification)**

- ▣ Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations
- ▣ New data is classified based on the training set

- **Unsupervised learning (clustering)**

- ▣ The class labels of training data is unknown
- ▣ Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

Prediction Problems: Classification vs. Numeric Prediction

□ Classification

- predicts categorical class labels (discrete or nominal)
- classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data

□ Numeric Prediction

- models continuous-valued functions, i.e., predicts unknown or missing values

Prediction Problems: Classification vs. Numeric Prediction

□ Classification

- predicts categorical class labels (discrete or nominal)
- classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data

□ Numeric Prediction

- models continuous-valued functions, i.e., predicts unknown or missing values

□ Typical applications

- Credit/loan approval:
- Medical diagnosis: if a tumor is cancerous or benign
- Fraud detection: if a transaction is fraudulent
- Web page categorization: which category it is

Classification—A Two-Step Process

(1) **Model construction:** describing a set of predetermined classes

- ▣ Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label** attribute
- ▣ The set of tuples used for model construction is **training set**
- ▣ Model: represented as classification rules, decision trees, or mathematical formulae

Classification—A Two-Step Process

(1) **Model construction:** describing a set of predetermined classes

- ▣ Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label** attribute
- ▣ The set of tuples used for model construction is **training set**
- ▣ Model: represented as classification rules, decision trees, or mathematical formulae

(2) **Model usage:** for classifying future or unknown objects

- ▣ Estimate accuracy of the model
 - The known label of test sample is compared with the classified result from the model
 - **Accuracy:** % of test set samples that are correctly classified by the model
 - Test set is independent of training set (otherwise overfitting)
- ▣ If the accuracy is acceptable, use the model to classify new data

Classification—A Two-Step Process

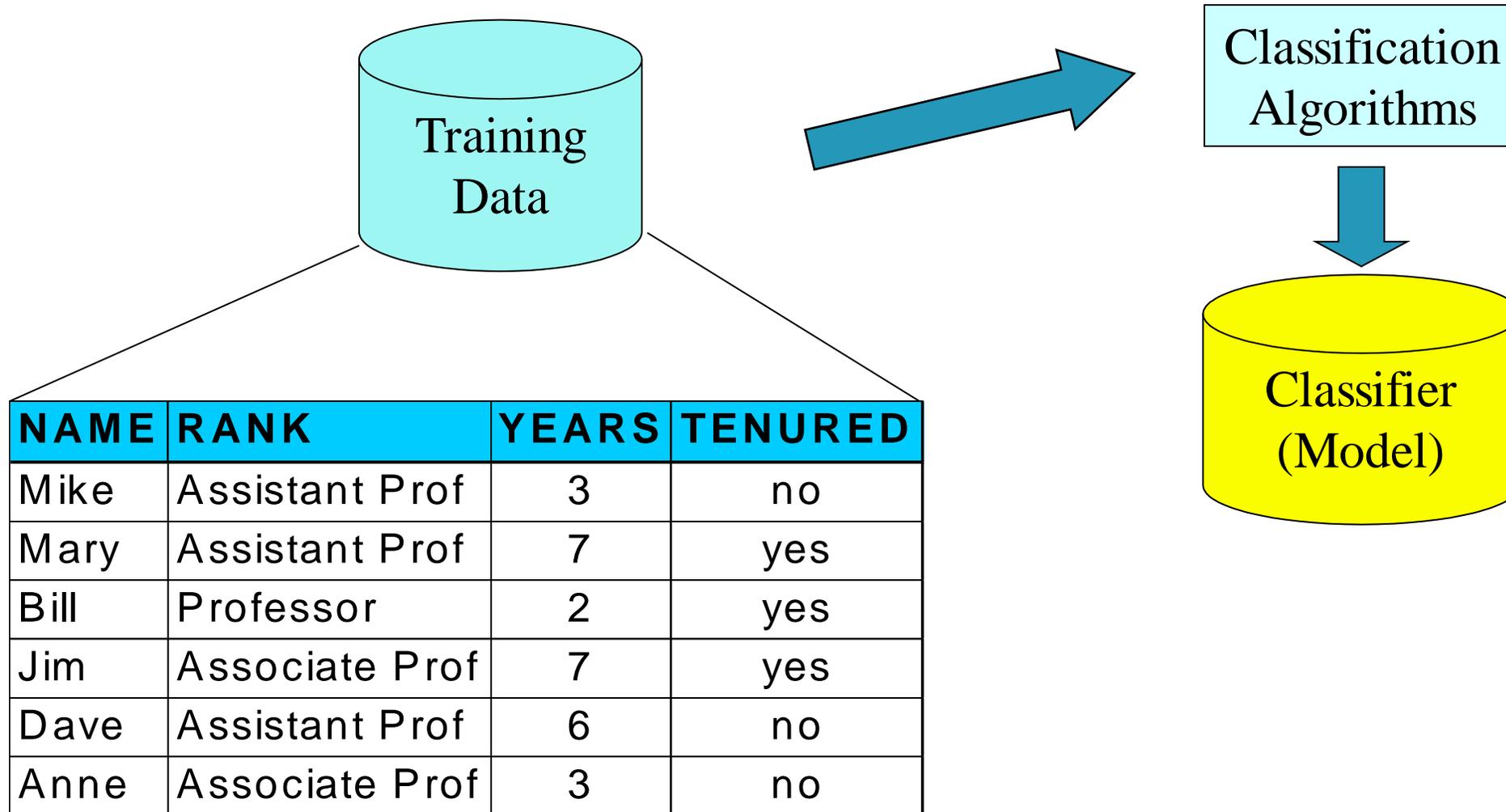
(1) **Model construction:** describing a set of predetermined classes

- ▣ Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label** attribute
- ▣ The set of tuples used for model construction is **training set**
- ▣ Model: represented as classification rules, decision trees, or mathematical formulae

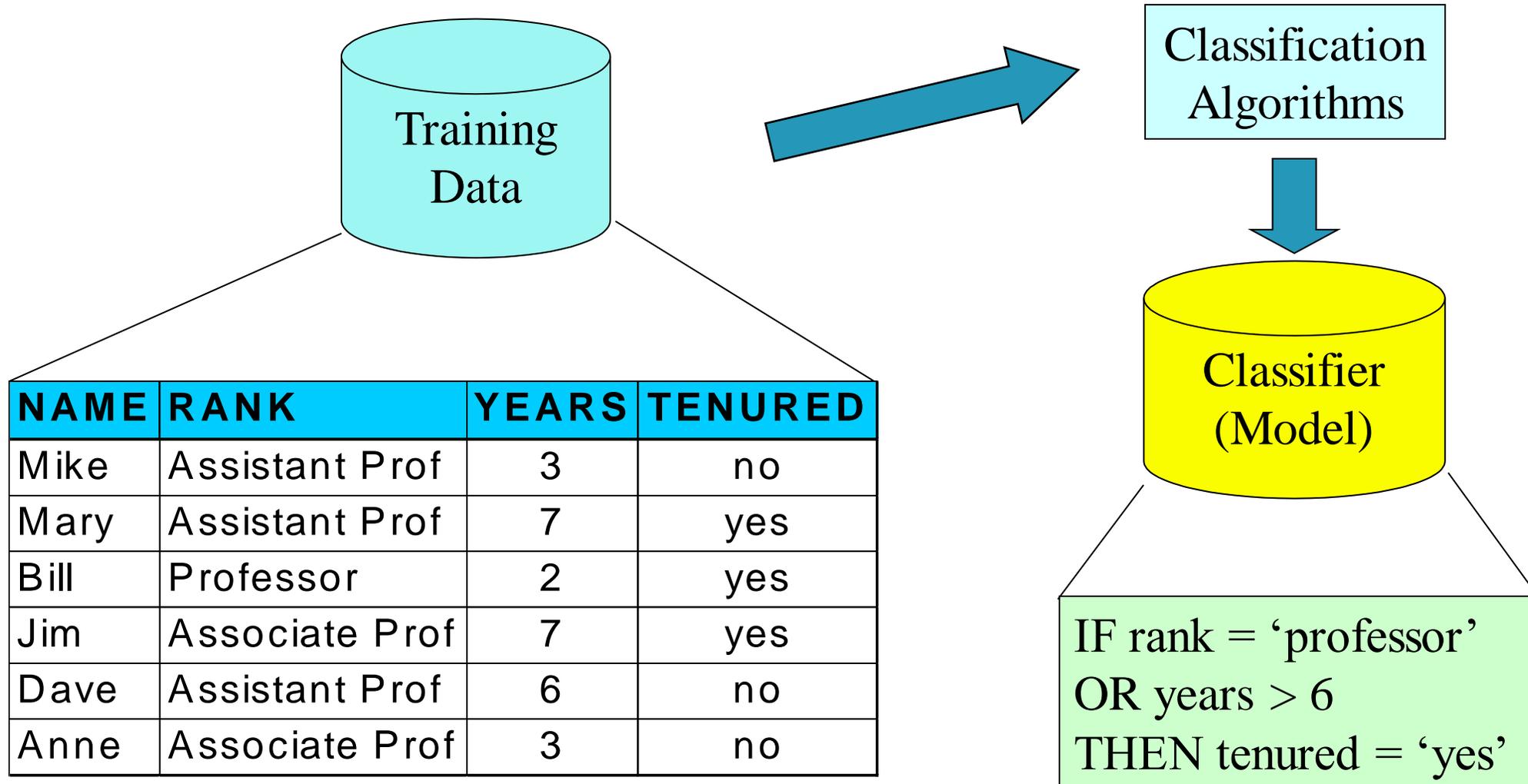
(2) **Model usage:** for classifying future or unknown objects

- ▣ Estimate accuracy of the model
 - The known label of test sample is compared with the classified result from the model
 - **Accuracy:** % of test set samples that are correctly classified by the model
 - Test set is independent of training set (otherwise overfitting)
- ▣ If the accuracy is acceptable, use the model to classify new data
- ▣ **Note:** If *the test set* is used to select/refine models, it is called **validation (test) set** or development test set

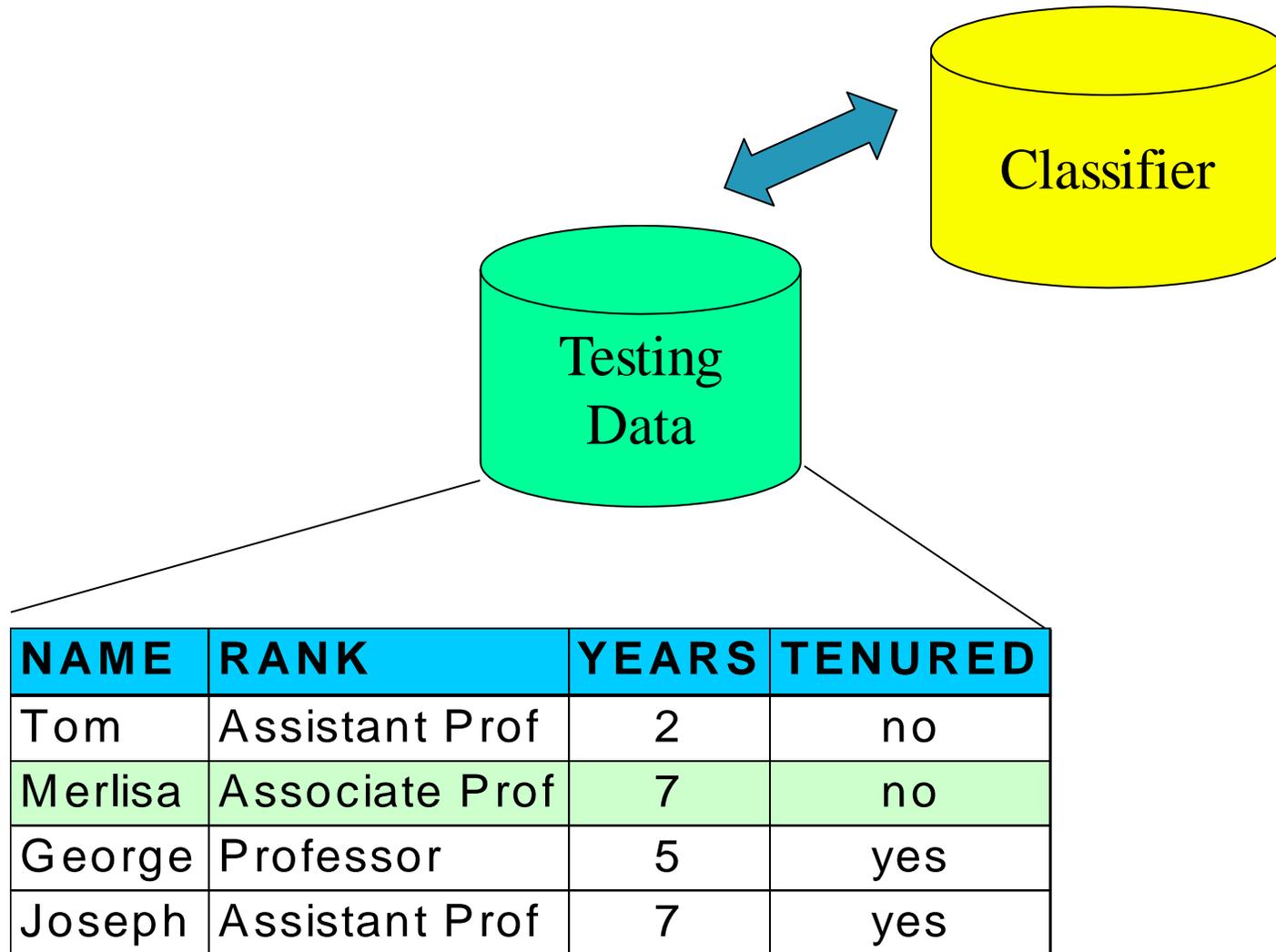
Step (1): Model Construction



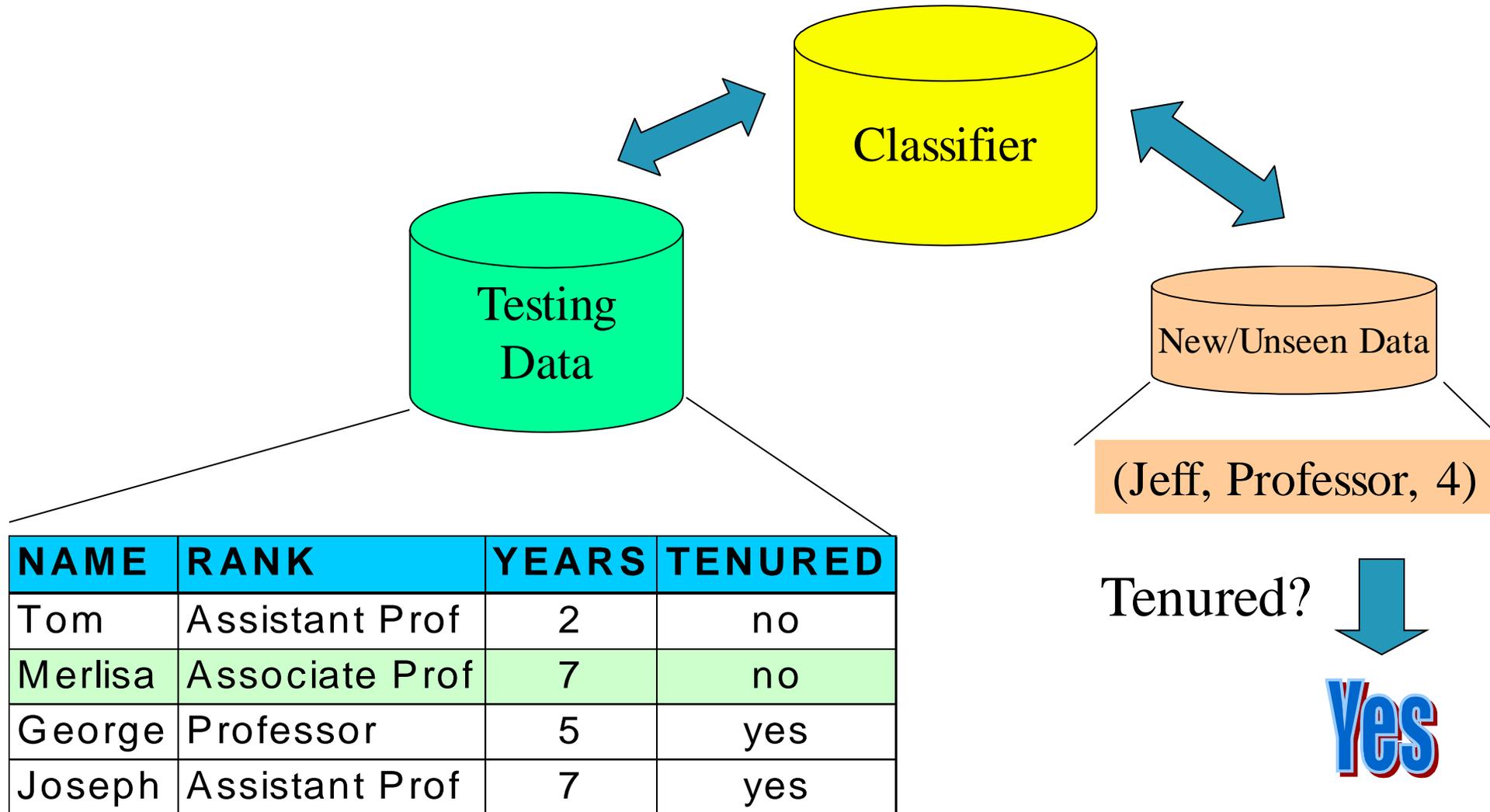
Step (1): Model Construction



Step (2): Using the Model in Prediction



Step (2): Using the Model in Prediction



Classification: Basic Concepts

- Classification: Basic Concepts
- Decision Tree Induction 
- Bayes Classification Methods
- Model Evaluation and Selection
- Techniques to Improve Classification Accuracy: Ensemble Methods
- Summary

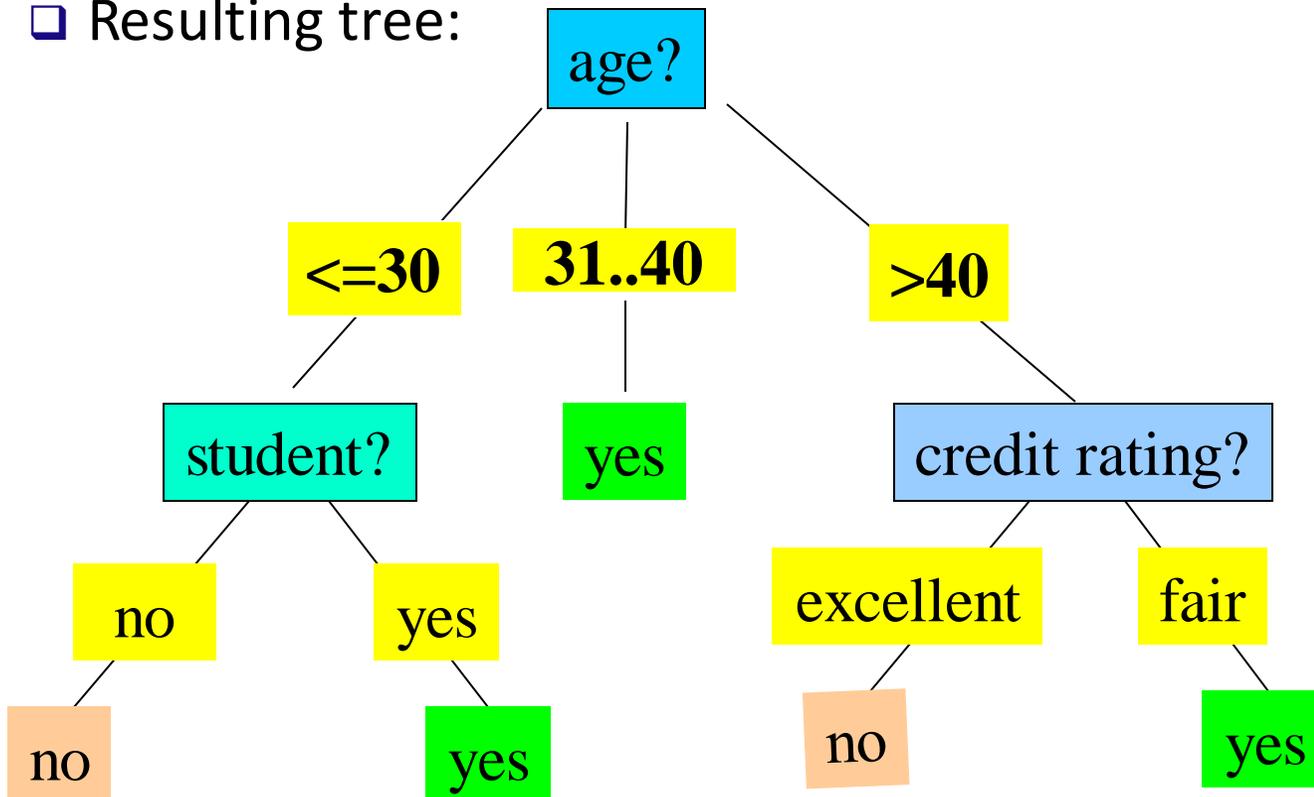
Decision Tree Induction: An Example

- ❑ Training data set: Buys_computer
- ❑ The data set follows an example of Quinlan's ID3 (Playing Tennis)

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Decision Tree Induction: An Example

- Training data set: Buys_computer
- The data set follows an example of Quinlan's ID3 (Playing Tennis)
- Resulting tree:



age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Algorithm for Decision Tree Induction

- Basic algorithm (a greedy algorithm)
 - ▣ Tree is constructed in a **top-down recursive divide-and-conquer manner**
 - ▣ At start, all the training examples are at the root
 - ▣ Attributes are categorical (if continuous-valued, they are discretized in advance)
 - ▣ Examples are partitioned recursively based on selected attributes
 - ▣ Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **information gain**)

Algorithm for Decision Tree Induction

- Basic algorithm (a greedy algorithm)
 - ▣ Tree is constructed in a **top-down recursive divide-and-conquer manner**
 - ▣ At start, all the training examples are at the root
 - ▣ Attributes are categorical (if continuous-valued, they are discretized in advance)
 - ▣ Examples are partitioned recursively based on selected attributes
 - ▣ Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **information gain**)
- Conditions for stopping partitioning
 - ▣ All samples for a given node belong to the same class
 - ▣ There are no remaining attributes for further partitioning—**majority voting** is employed for classifying the leaf
 - ▣ There are no samples left

Algorithm Outline

- Split (node, {data tuples})
 - ▣ $A \leq$ the best attribute for splitting the {data tuples}
 - ▣ Decision attribute for this node $\leq A$
 - ▣ For each value of A , create new child node
 - ▣ For each child node / subset:
 - If one of the stopping conditions is satisfied: STOP
 - Else: Split (child_node, {subset})

ID3 algorithm: how it works

https://www.youtube.com/watch?v=_XhOdSLIE5c

Algorithm Outline

- Split (node, {data tuples})
 - ▣ $A \leq$ the **best attribute** for splitting the {data tuples}
 - ▣ Decision attribute for this node $\leq A$
 - ▣ For each value of A , create new child node
 - ▣ For each child node / subset:
 - If one of the stopping conditions is satisfied: STOP
 - Else: Split (child_node, {subset})

ID3 algorithm: how it works

https://www.youtube.com/watch?v=_XhOdSLIE5c

Brief Review of Entropy

□ Entropy (Information Theory)

▣ A measure of uncertainty associated with a random number

▣ Calculation: For a discrete random variable Y taking m distinct values $\{y_1, y_2, \dots, y_m\}$

$$H(Y) = - \sum_{i=1}^m p_i \log(p_i) \quad \text{where } p_i = P(Y = y_i)$$

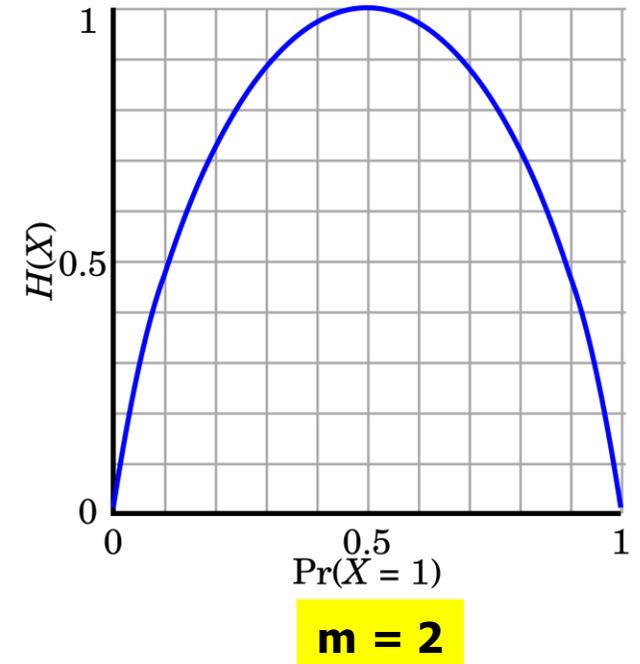
▣ Interpretation

■ Higher entropy \rightarrow higher uncertainty

■ Lower entropy \rightarrow lower uncertainty

□ Conditional entropy

$$H(Y|X) = \sum_x p(x) H(Y|X = x)$$



Attribute Selection Measure: Information Gain (ID3/C4.5)

- ❑ Select the attribute with the highest information gain
- ❑ Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$
- ❑ Expected information (entropy) needed to classify a tuple in D :

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$

- ❑ Information needed (after using A to split D into v partitions) to classify D :

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

- ❑ Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

How to select the first attribute?

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$$Inf\alpha(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$$Inf\alpha(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Look at "age":

age	p_i	n_i	$I(p_i, n_i)$
<=30	2	3	0.971
31...40	4	0	0
>40	3	2	0.971

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Look at "age":

age	p_i	n_i	$I(p_i, n_i)$
<=30	2	3	0.971
31...40	4	0	0
>40	3	2	0.971

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

Look at "age":

age	p_i	n_i	$I(p_i, n_i)$
≤ 30	2	3	0.971
31...40	4	0	0
> 40	3	2	0.971

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

$\frac{5}{14} I(2,3)$ means "age ≤ 30 " has 5 out of 14 samples, with 2 yes'es and 3 no's.

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

Similarly,

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$

Recursive Procedure

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

1. After selecting age at the root node, we will create three child nodes.
2. One child node is associated with red data tuples.
3. How to continue for this child node?

Recursive Procedure

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

1. After selecting age at the root node, we will create three child nodes.

2. One child node is associated with red data tuples.

3. How to continue for this child node?

Now, you will make $D = \{\text{red data tuples}\}$

and then select the best attribute to further split

D.

A recursive procedure.

Computing Information-Gain for Continuous-Valued Attributes

- Let attribute A be a continuous-valued attribute

Computing Information-Gain for Continuous-Valued Attributes

- Let attribute A be a continuous-valued attribute
- Must determine the **best split point** for A (or, discretization)

\Rightarrow

D_1 is the set of tuples in D satisfying $A \leq \text{split-point}$, and D_2 is the set of tuples in D satisfying $A > \text{split-point}$

Computing Information-Gain for Continuous-Valued Attributes

- Let attribute A be a continuous-valued attribute
- Must determine the **best split point** for A
 - ▣ Sort the value A in increasing order
 - ▣ Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*
 - $(a_i + a_{i+1})/2$ is the midpoint between the values of a_i and a_{i+1}
 - ▣ The point with the **minimum expected information requirement for A (i.e., $\text{Info}_A(D)$)** is selected as the split-point for A

Computing Information-Gain for Continuous-Valued Attributes

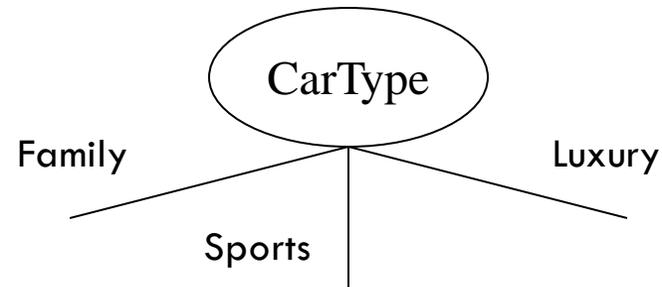
- Let attribute A be a continuous-valued attribute
- Must determine the **best split point** for A
 - Sort the value A in increasing order
 - Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*
 - $(a_i + a_{i+1})/2$ is the midpoint between the values of a_i and a_{i+1}
 - The point with the **minimum expected information requirement for A (i.e., $\text{Info}_A(D)$)** is selected as the split-point for A
- Split:
 - D_1 is the set of tuples in D satisfying $A \leq \text{split-point}$, and D_2 is the set of tuples in D satisfying $A > \text{split-point}$

How to Select Test Attribute?

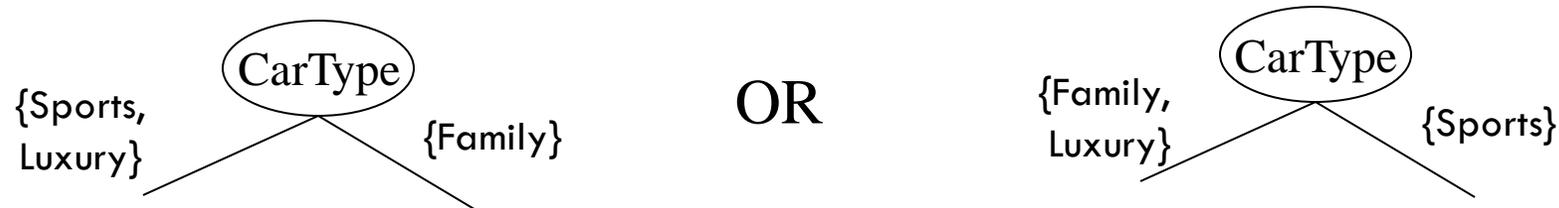
- Depends on attribute types
 - ▣ Nominal
 - ▣ Ordinal
 - ▣ Continuous
- Depends on number of ways to split
 - ▣ 2-way split
 - ▣ Multi-way split

Splitting Based on Nominal Attributes

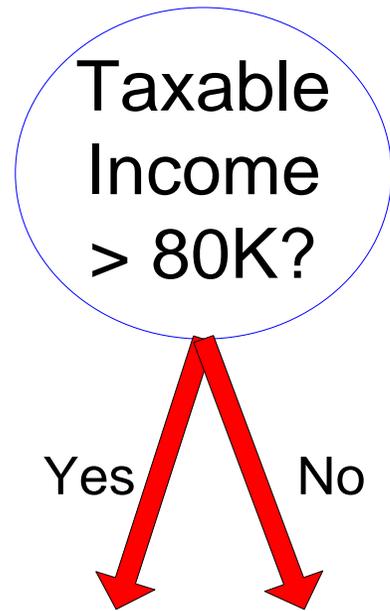
- **Multi-way split:** Use as many partitions as distinct values.



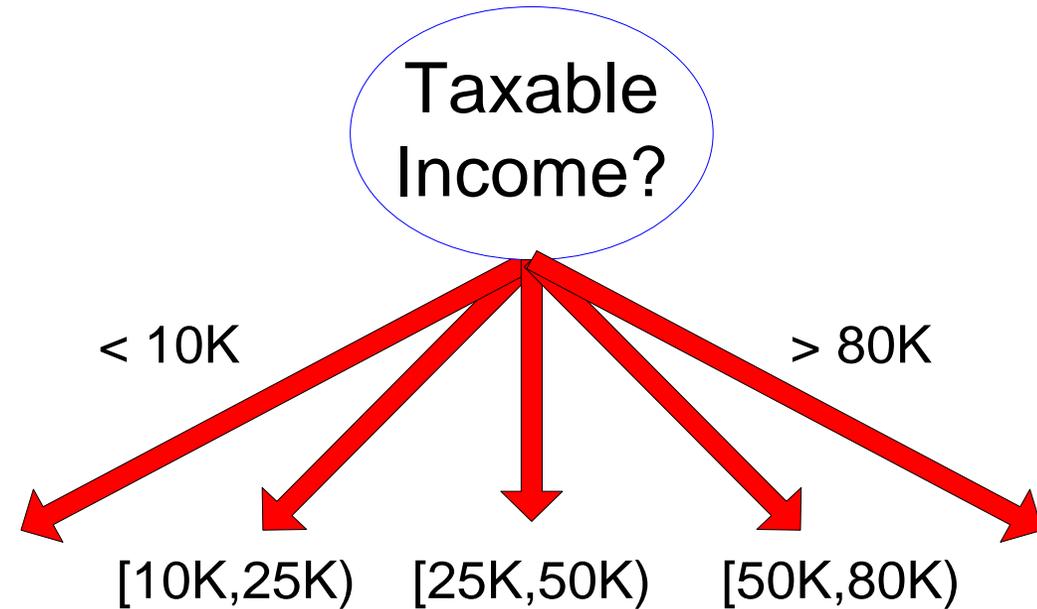
- **Binary split:** Divides values into two subsets. **Need to find optimal partitioning.**



Splitting Based on Continuous Attributes



(i) Binary split

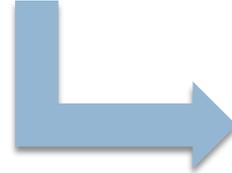


(ii) Multi-way split

How to Determine the Best Split

- Greedy approach:

- Nodes with **homogeneous** class distribution are preferred



Ideally, data tuples at that node belong to the same class.

C0: 5
C1: 5

Non-homogeneous,
High degree of impurity

C0: 9
C1: 1

Homogeneous,
Low degree of impurity

Rethink about Decision Tree Classification

- 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

Measures of Node Impurity

- Entropy:
$$H(Y) = - \sum_{i=1}^m p_i \log(p_i) \quad \text{where } p_i = P(Y = y_i)$$
 - ▣ Higher entropy => higher uncertainty, higher node impurity
 - ▣ Why entropy is used in information gain
- Gini Index
- Misclassification error

Potential Problems with Information Gain

- Information gain measure is **biased** towards **attributes with a large number of values (or, a large number of small partitions)**

Potential Problems with Information Gain

- Information gain measure is biased towards attributes with a large number of values (or, a large number of small partitions)
 - ▣ Consider an attribute that acts as a **unique identifier**, such as student_ID

Potential Problems with Information Gain

- Information gain measure is biased towards attributes with a large number of values (or, a large number of small partitions)
 - ▣ Consider an attribute that acts as a **unique identifier**, such as student_ID
 - ▣ A split on Student_ID will result in many partitions, each one containing just one tuple.

Potential Problems with Information Gain

- Information gain measure is biased towards attributes with a large number of values (or, a large number of small partitions)
 - ▣ Consider an attribute that acts as a **unique identifier**, such as student_ID
 - ▣ A split on Student_ID will result in many partitions, each one containing just one tuple.
 - ▣ Information required to classify data set D after this partitioning is:

?

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

Gain Ratio for Attribute Selection (C4.5)

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

Gain Ratio for Attribute Selection (C4.5)

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$\text{SplitInfo}_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

- The entropy of the partitioning, or the potential information generated by splitting D into v partitions.

Gain Ratio for Attribute Selection (C4.5)

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$\text{SplitInfo}_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

- The entropy of the partitioning, or the potential information generated by splitting D into v partitions.
- **GainRatio(A) = Gain(A)/SplitInfo(A)** (normalizing Information Gain)

Gain Ratio for Attribute Selection (C4.5)

- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$\text{SplitInfo}_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

- GainRatio(A) = Gain(A)/SplitInfo(A)
- Ex.

$$\text{Gain}(\text{income}) = 0.029 \quad (\text{from last class, slide 27})$$

$$\text{SplitInfo}_{\text{income}}(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \times \log_2\left(\frac{4}{14}\right) = 1.557$$

- gain_ratio(income) = 0.029/1.557 = 0.019
- The attribute with the **maximum gain ratio** is selected as the splitting attribute

Another Attribute Selection Measure: Gini Index (CART, IBM IntelligentMiner)

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^n P_j^2$$

where p_j is the relative frequency of class j in D

Gini index measures the impurity of D . The larger, the more impure.

Another Attribute Selection Measure: Gini Index (CART, IBM IntelligentMiner)

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^n P_j^2$$

where p_j is the relative frequency of class j in D

Ex. D has 9 tuples in `buys_computer = "yes"` and 5 in `"no"`

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

Gini Index (CART, IBM IntelligentMiner)

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^n p_j^2, \text{ where } p_j \text{ is the relative frequency of class } j \text{ in } D$$

- If a data set D is split on A into two subsets D_1 and D_2 , the gini index after the split is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

Gini Index (CART, IBM IntelligentMiner)

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^n p_j^2, \text{ where } p_j \text{ is the relative frequency of class } j \text{ in } D$$

- If a data set D is split on A into two subsets D_1 and D_2 , the gini index $gini(D)$ is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- Reduction in impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

Gini Index (CART, IBM IntelligentMiner)

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^n p_j^2, \text{ where } p_j \text{ is the relative frequency of class } j \text{ in } D$$

- If a data set D is split on A into two subsets D_1 and D_2 , the *gini index after the split* is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- Reduction in impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

- The attribute provides the smallest $gini_A(D)$ (or, **the largest reduction in impurity**) is chosen to split the node.

Gini Index (CART, IBM IntelligentMiner)

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^n p_j^2, \text{ where } p_j \text{ is the relative frequency of class } j \text{ in } D$$

- If a data set D is split on A into two subsets D_1 and D_2 , the *gini index after the split* is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

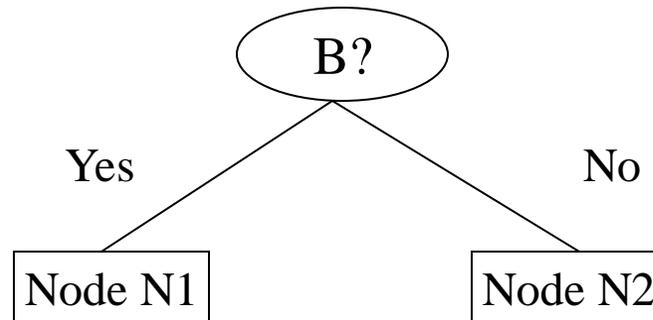
- Reduction in impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

- The attribute provides the smallest $gini_A(D)$ (or, **the largest reduction in impurity**) is chosen to split the node.

Binary Attributes: Computing Gini Index

- Splits into two partitions
- Effect of weighing partitions:
 - Larger and Purer Partitions are sought for.



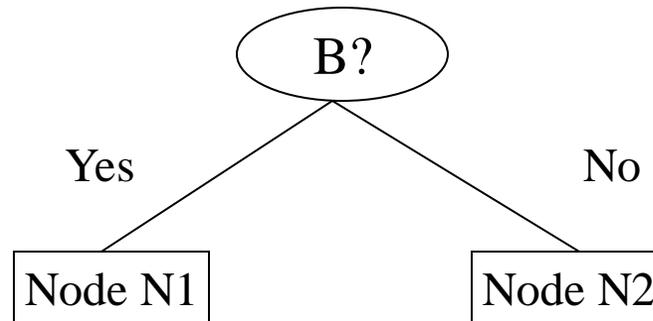
$$gini(D) = 1 - \sum_{j=1}^n p_j^2$$

	Parent
C1	6
C2	6
Gini = ?	

Binary Attributes: Computing Gini Index

- Splits into two partitions
- Effect of weighing partitions:
 - Larger and Purer Partitions are sought for.

$$gini(D) = 1 - \sum_{j=1}^n p_j^2$$



	Parent
C1	6
C2	6
Gini = 0.500	

$$\begin{aligned} \text{Gini}(N1) &= 1 - (5/6)^2 - (2/6)^2 \\ &= 0.194 \end{aligned}$$

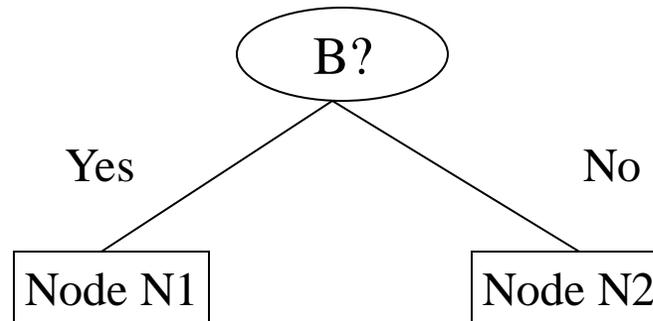
$$\begin{aligned} \text{Gini}(N2) &= 1 - (1/6)^2 - (4/6)^2 \\ &= 0.528 \end{aligned}$$

	N1	N2
C1	5	1
C2	2	4
Gini=?		

Binary Attributes: Computing Gini Index

- Splits into two partitions
- Effect of weighing partitions:
 - Larger and Purer Partitions are sought for.

$$gini(D) = 1 - \sum_{j=1}^n p_j^2$$



	Parent
C1	6
C2	6
Gini = ?	

$$\begin{aligned} \text{Gini}(N1) &= 1 - (5/6)^2 - (2/6)^2 \\ &= 0.194 \end{aligned}$$

$$\begin{aligned} \text{Gini}(N2) &= 1 - (1/6)^2 - (4/6)^2 \\ &= 0.528 \end{aligned}$$

	N1	N2
C1	5	1
C2	2	4
Gini=0.333		

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



weighting

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)

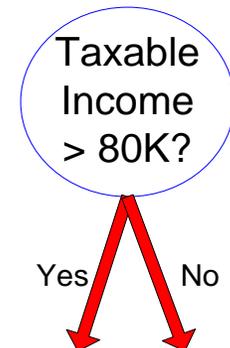
	CarType	
	{Sports, Luxury}	{Family}
C1	3	1
C2	2	4
Gini	0.400	

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

Continuous Attributes: Computing Gini Index

- Use Binary Decisions based on one value
- Several Choices for the splitting value
 - Number of possible splitting values = Number of distinct values - 1
 - Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*
 - $(a_i + a_{i+1})/2$ is the midpoint between the values of a_i and a_{i+1}
- Each splitting value has a count matrix associated with it
 - Class counts in each of the partitions, $A < v$ and $A \geq v$
- Simple method to choose best v
 - For each v , scan the database to gather count matrix and compute its Gini index
 - Computationally Inefficient! Repetition of work.

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



Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
 - ▣ Sort the attribute on values
 - ▣ Linearly 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											
Sorted Values	60	70	75	85	90	95	100	120	125	220	
Split Positions		65	72	80	87	92	97	110	122	172	
		<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >	
Yes		0 3	0 3	0 3	1 2	2 1	3 0	3 0	3 0	3 0	
No		1 6	2 5	3 4	3 4	3 4	3 4	4 3	5 2	6 1	
Gini		0.400	0.375	0.343	0.417	0.400	<u>0.300</u>	0.343	0.375	0.400	

Attribute Selection: Complexity Note

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Look at "age":

age	p_i	n_i
<=30	2	3
31...40	4	0
>40	3	2

When selecting an attribute, scan the data set at each node to gather count matrix



Another Impurity Measure: Misclassification Error

- Classification error at a node t :

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

- $P(i | t)$ means the relative frequency of class i at node t .
- Measures misclassification error made by a node.
 - Maximum $(1 - 1/n_c)$ when records are equally distributed among all classes, implying most impurity
 - Minimum (0.0) when all records belong to one class, implying least impurity

Examples for Misclassification Error

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

C1	0
C2	6

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

$$Error = 1 - \max(0, 1) = 1 - 1 = 0$$

C1	1
C2	5

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

$$Error = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

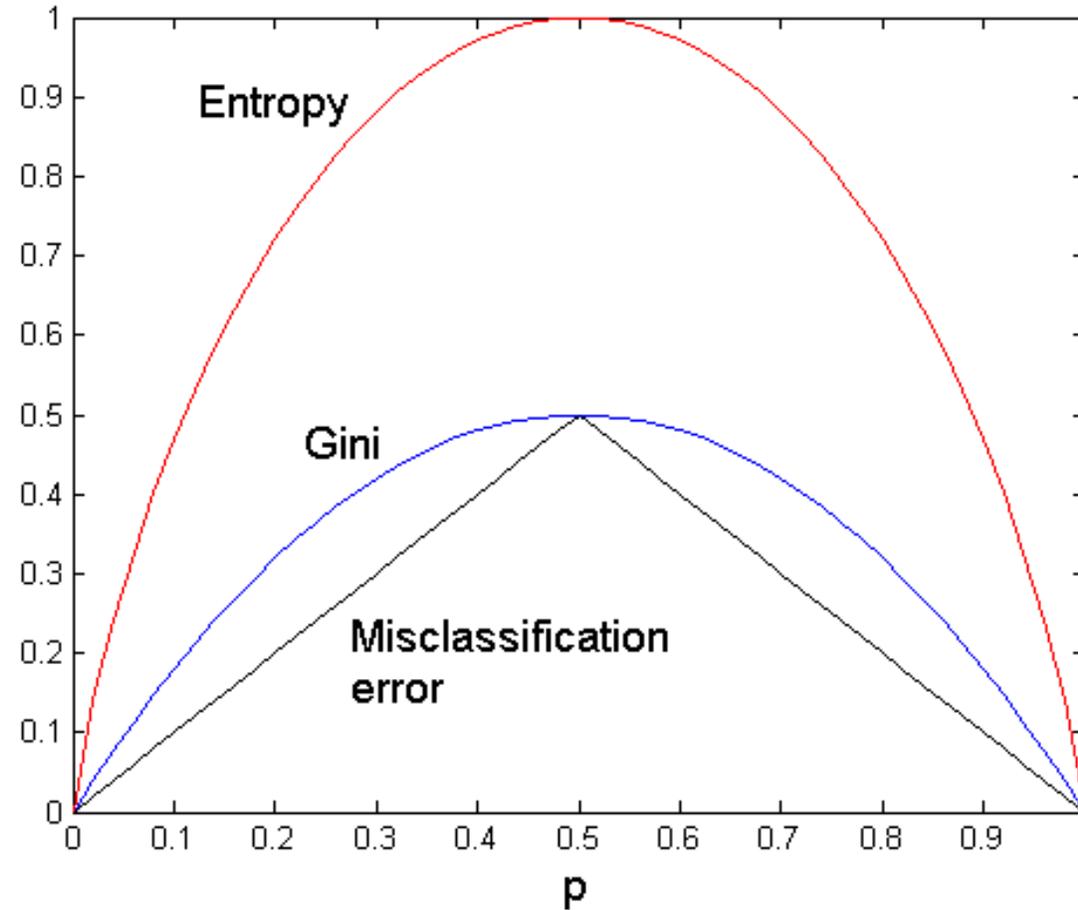
C1	2
C2	4

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

$$Error = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$

Comparison among Impurity Measure

For a 2-class problem:



Comparing Attribute Selection Measures

- The three measures, in general, return good results but
 - ▣ **Information gain:**
 - biased towards multivalued attributes
 - ▣ **Gain ratio:**
 - tends to prefer unbalanced splits in which one partition is much smaller than the others
 - ▣ **Gini index:**
 - biased to multivalued attributes
 - has difficulty when # of classes is large
 - tends to favor tests that result in equal-sized partitions and purity in both partitions

Other Attribute Selection Measures

- CHAID: a popular decision tree algorithm, measure based on χ^2 test for independence
- C-SEP: performs better than info. gain and gini index in certain cases
- G-statistic: has a close approximation to χ^2 distribution
- MDL (Minimal Description Length) principle (i.e., the simplest solution is preferred):
 - ▣ The best tree as the one that requires the fewest # of bits to both (1) encode the tree, and (2) encode the exceptions to the tree
- Multivariate splits (partition based on multiple variable combinations)
 - ▣ CART: finds multivariate splits based on a linear comb. of attrs.
- **Which attribute selection measure is the best?**
 - ▣ Most give good results, none is significantly superior than others

Decision Tree Based Classification

- Advantages:
 - ▣ Inexpensive to construct
 - ▣ Extremely fast at classifying unknown records
 - ▣ Easy to interpret for small-sized trees
 - ▣ Accuracy is comparable to other classification techniques for many simple data sets

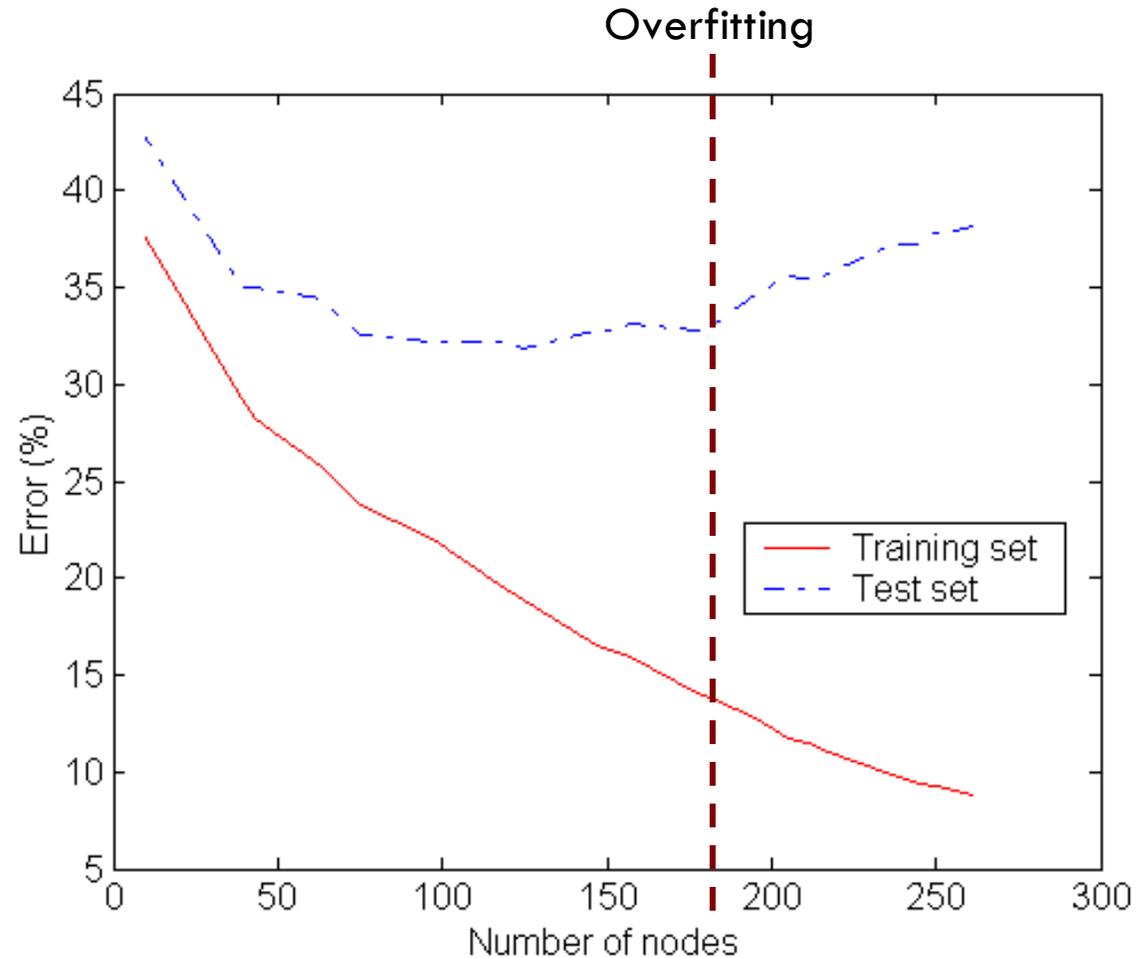
Example: C4.5

- Simple depth-first construction.
- Uses Information Gain
- **Sorts Continuous Attributes at each node.**
- Needs entire data to fit in memory.
- Unsuitable for Large Datasets.
 - ▣ Needs out-of-core sorting.
- You can download the software online, e.g.,
<http://www2.cs.uregina.ca/~dbd/cs831/notes/ml/dtrees/c4.5/tutorial.html>

Practical Issues of Classification

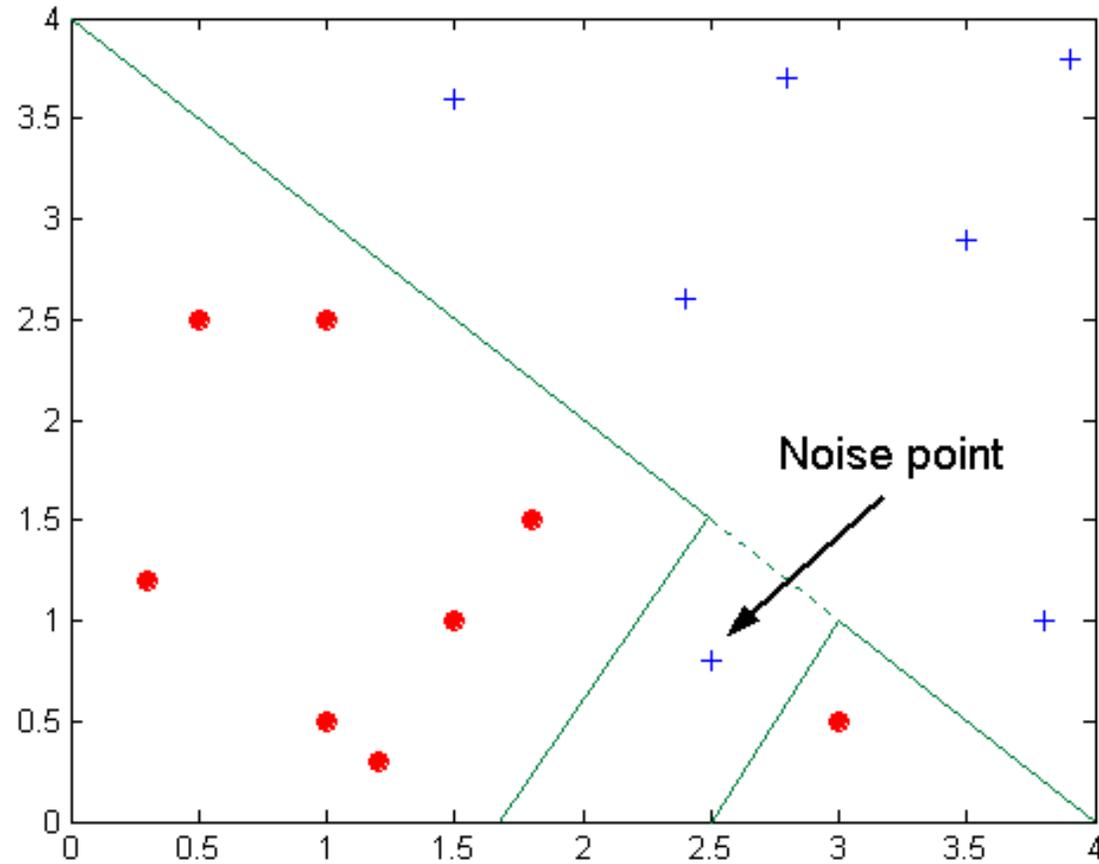
- Underfitting and Overfitting
- Missing Values
- Costs of Classification

Underfitting and Overfitting



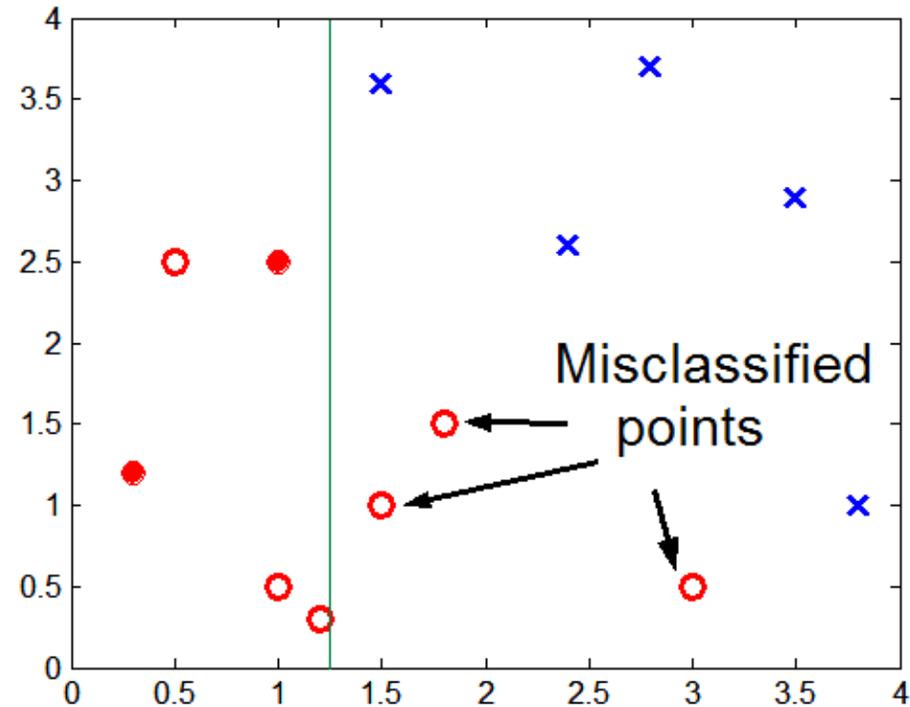
Underfitting: when model is too simple, both training and test errors are large

Overfitting due to Noise



Decision boundary is distorted by noise point

Overfitting due to Insufficient Examples



Lack of data points in the lower half of the diagram makes it difficult to predict correctly the class labels of that region

- Insufficient number of training records in the region causes the decision tree to predict the test examples using other training records that are irrelevant to the classification task

Notes on Overfitting

- Overfitting results in decision trees that are more complex than necessary
- Training error no longer provides a good estimate of how well the tree will perform on previously unseen records
- Need new ways for estimating errors