



# Classification

---

CS145  
Fall 2015

# Classification vs. Prediction

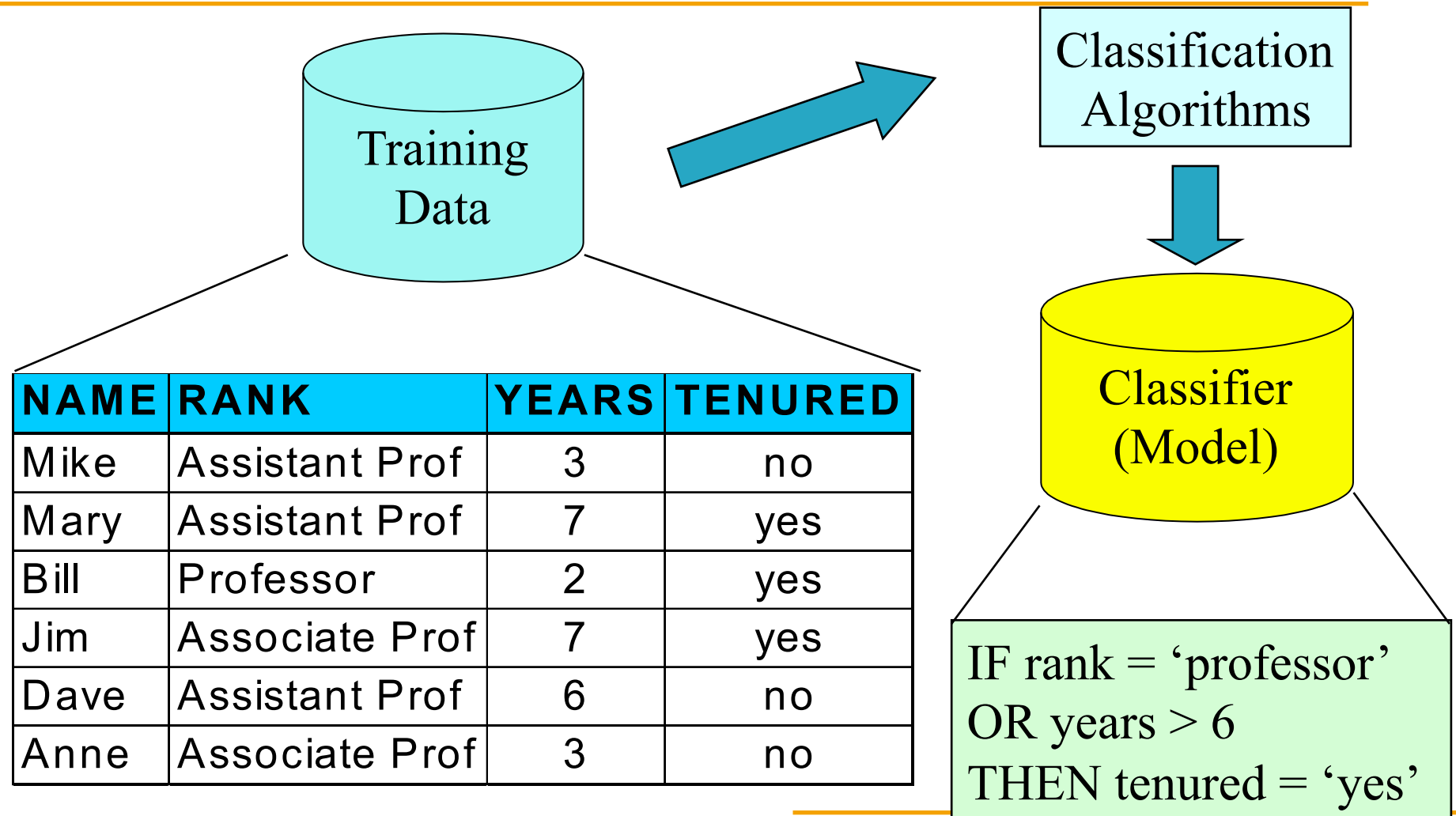
---

- ▶ **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
- ▶ **Typical Applications**
  - ▶ credit approval
  - ▶ target marketing
  - ▶ medical diagnosis
  - ▶ treatment effectiveness analysis

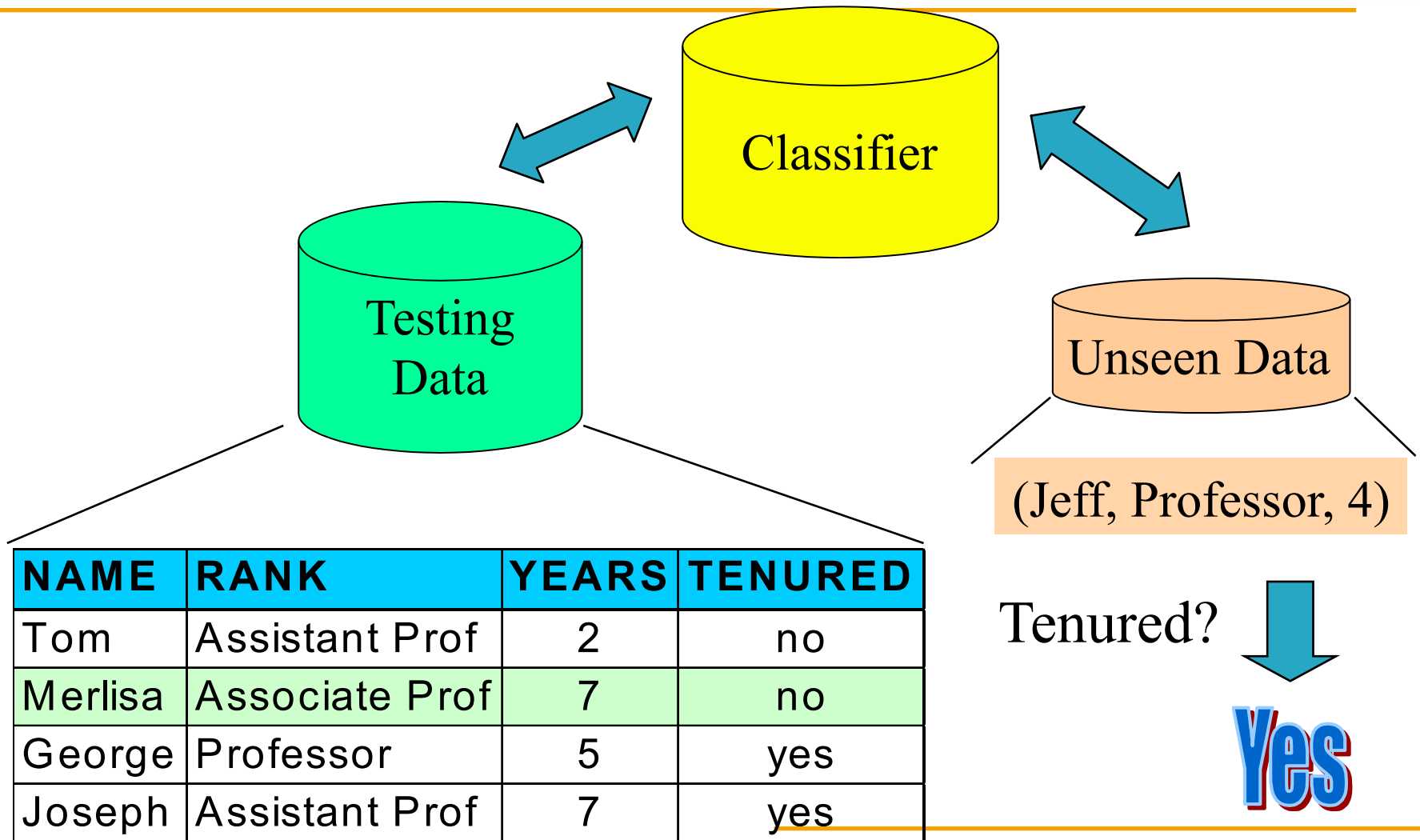
# Classification—A Two-Step Process

- ▶ **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**
  - ▶ The model is represented as classification rules, decision trees, or mathematical formulae
- ▶ **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 rate is the percentage of test set samples that are correctly classified by the model
    - ▶ Test set is independent of training set
    - ▶ If the accuracy is acceptable, use the model to classify data tuples whose class labels are not known

# Classification Process (1): Model Construction



# Classification Process (2): Use the Model in Prediction



# 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

# Major Classification Models

---

- ▶ Classification by decision tree induction
- ▶ Bayesian Classification
- ▶ Neural Networks
- ▶ Support Vector Machines (SVM)
- ▶ Classification Based on Associations
- ▶ Other Classification Methods
  - ▶ KNN
  - ▶ Boosting
  - ▶ Bagging
  - ▶ ...

# Evaluating Classification Methods

---

- ▶ Predictive accuracy
- ▶ Speed and scalability
  - ▶ time to construct the model
  - ▶ time to use the model
- ▶ Robustness
  - ▶ handling noise and missing values
- ▶ Scalability
  - ▶ efficiency in disk-resident databases
- ▶ Interpretability:
  - ▶ understanding and insight provided by the model
- ▶ Goodness of rules
  - ▶ decision tree size
  - ▶ compactness of classification rules



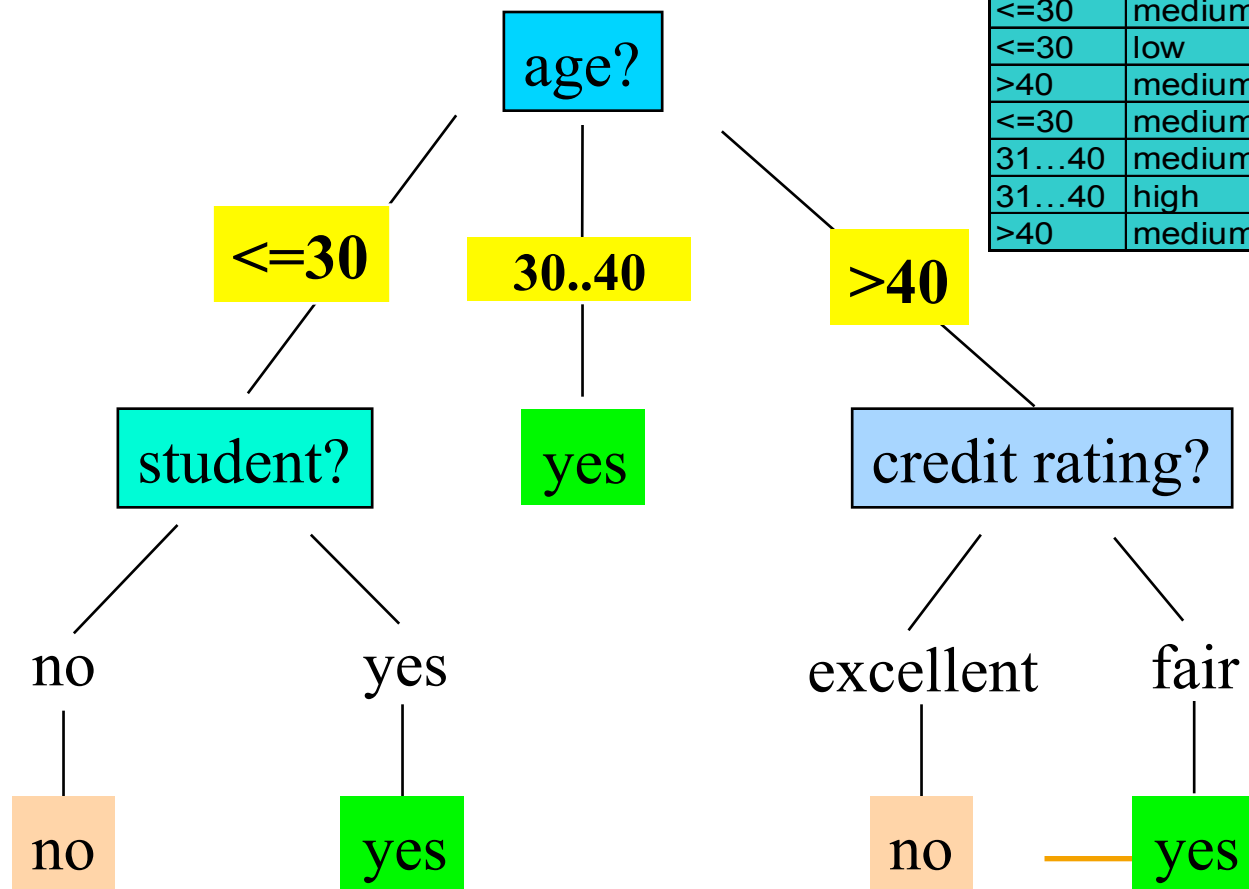
# Decision Tree

## Training Dataset

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

# Output: A Decision Tree for “*buys\_computer*”

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**)
- ▶ 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

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

- Select the attribute with the highest information gain
- S contains  $s_i$  tuples of class  $C_i$  for  $i = \{1, \dots, m\}$
- **information** measures info required to classify any arbitrary tuple

$$I(s_1, s_2, \dots, s_m) = - \sum_{i=1}^m \frac{s_i}{s} \log_2 \frac{s_i}{s}$$

- **entropy** of attribute  $A$  with values  $\{a_1, a_2, \dots, a_v\}$

$$E(A) = \sum_{j=1}^v \frac{s_{1j} + \dots + s_{mj}}{s} I(s_{1j}, \dots, s_{mj})$$

- **information gained** by branching on attribute  $A$

$$\text{Gain}(A) = I(s_1, s_2, \dots, s_m) - E(A)$$

# Attribute Selection by Information Gain Computation

■ Class P: buys\_computer = “yes”

■ Class N: buys\_computer = “no”

■  $I(p, n) = I(9, 5) = 0.940$

■ Compute the entropy for *age*:

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

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

$\frac{5}{14} I(2,3)$  means “age  $\leq 30$ ” has 5 out of 14 samples, with 2 yes’es and 3 no’s. Hence

$$\text{Gain}(\text{age}) = I(p, n) - E(\text{age}) = 0.246$$

Similarly,

$$\text{Gain}(\text{income}) = 0.029$$

$$\text{Gain}(\text{student}) = 0.151$$

$$\text{Gain}(\text{credit\_rating}) = 0.048$$

age	income	student	credit_rating	buys_computer
$\leq 30$	high	no	fair	no
$\leq 30$	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
$\leq 30$	medium	no	fair	no
$\leq 30$	low	yes	fair	yes
$> 40$	medium	yes	fair	yes
$\leq 30$	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
$> 40$	medium	no	excellent	no

# Natural Bias in The Information Gain Measure

---

- ▶ Favor attributes with many values
- ▶ An extreme example
  - ▶ Attribute “income” might have the highest information gain
  - ▶ A very broad decision tree of depth one
  - ▶ Inapplicable to any future data

# Alternative Measures

- ▶ Gain ratio: penalize attributes like income by incorporating split information
  - ▶  $SplitInformation(S, A) \equiv - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$ 
    - ▶ Split information is sensitive to how broadly and uniformly the attribute splits the data
  - ▶  $GainRatio(S, A) \equiv \frac{Gain(S, A)}{SplitInformation(S, A)}$
- ▶ Gain ratio can be undefined or very large
  - ▶ Only test attributes with above average Gain

# Other Attribute Selection Measures

---

- ▶ **Gini index** (CART, IBM IntelligentMiner)
  - ▶ All attributes are assumed continuous-valued
  - ▶ Assume there exist several possible split values for each attribute
  - ▶ May need other tools, such as clustering, to get the possible split values
  - ▶ Can be modified for categorical attributes



# Gini Index (IBM IntelligentMiner)

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

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

where  $p_j$  is the relative frequency of class  $j$  in  $T$ .

- ▶ If a data set  $T$  is split into two subsets  $T_1$  and  $T_2$  with sizes  $N_1$  and  $N_2$  respectively, the gini index of the split data contains examples from  $n$  classes, the gini index  $gini(T)$  is defined as

$$gini_{split}(T) = \frac{N_1}{N} gini(T_1) + \frac{N_2}{N} gini(T_2)$$

- ▶ The attribute provides the smallest  $gini_{split}(T)$  is chosen to split the node (*need to enumerate all possible splitting points for each attribute*).

# Extracting Classification Rules from Trees

- ▶ Represent the knowledge in the form of **IF-THEN** rules
- ▶ One rule is created for each path from the root to a leaf
- ▶ Each attribute-value pair along a path forms a conjunction
- ▶ The leaf node holds the class prediction
- ▶ Rules are easier for humans to understand
- ▶ Example

IF *age* = “≤30” AND *student* = “no” THEN *buys\_computer* = “no”

IF *age* = “≤30” AND *student* = “yes” THEN *buys\_computer* = “yes”

IF *age* = “31...40” THEN *buys\_computer* = “yes”

IF *age* = “>40” AND *credit\_rating* = “excellent” THEN *buys\_computer* = “yes”

IF *age* = “>40” AND *credit\_rating* = “fair” THEN *buys\_computer* = “no”

# Avoid Overfitting in Classification

---

- ▶ Overfitting: An induced tree may overfit the training data
  - ▶ Too many branches, some may reflect anomalies due to noise or outliers
  - ▶ Poor accuracy for unseen samples
- ▶ Two approaches to avoid overfitting
  - ▶ Prepruning: Halt tree construction early—do not split a node if this would result in the goodness measure falling below a threshold
    - ▶ Difficult to choose an appropriate threshold
  - ▶ Postpruning: Remove branches from a “fully grown” tree—get a sequence of progressively pruned trees
    - ▶ Use a set of data different from the training data to decide which is the “best pruned tree”

# Approaches to Determine the Final Tree Size

---

- ▶ Separate training (2/3) and testing (1/3) sets
- ▶ Use cross validation, e.g., 10-fold cross validation
- ▶ Use all the data for training
  - ▶ but apply a **statistical test** (e.g., chi-square) to estimate whether expanding or pruning a node may improve the entire distribution
- ▶ Use minimum description length (MDL) principle
  - ▶ halting growth of the tree when the encoding is minimized

# Minimum Description Length

---

- ▶ The ideal MDL select the model with the shortest effective description that minimizes the sum of
  - ▶ The length, in bits, of an effective description of the model; and
  - ▶ The length, in bits, of an effective description of the data when encoded with help of the model

$$H_0 = \min_{H \in \mathcal{H}} \{K(D|H) + K(H)\}$$