

## Outline

- ▶ Overview of classification
- ▶ Basic concept in decision tree learning
  - ▶ Data representation in decision tree learning
  - ▶ What is a decision tree?
- ▶ How to learn a decision tree from data
  - ▶ Basic decision tree learning algorithm
  - ▶ How to select best attribute
    - ▶ Information gain
    - ▶ Gain ratio
    - ▶ Gini index
  - ▶ Pruning decision tree
    - ▶ Pre-pruning
    - ▶ Post-pruning
- ▶ Other issues involved in decision tree learning

24

## Bias in the Information Gain Measure

- ▶ Favor unfairly attributes with large numbers of distinct values at the expense of those with few.
  - ▶ E.g., attribute Date: poor predictor, but has the highest gain because it alone perfectly predicts the target attribute over the training data

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	Rainy	Mild	High	Weak	Yes
D5	Rainy	Cool	Normal	Weak	Yes
D6	Rainy	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rainy	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rainy	Mild	High	Strong	No

25

## Gain Ratio

- ▶ Proposed by Quinlan in 1986 (used in C4.5)
- ▶ Idea:
  - ▶ penalizes attributes with many distinct values by dividing information gain by attribute information (entropy of data with respect to the values of attribute):

$$SplitInformation(S, A) = - \sum_{v_i \in Values(A)} \frac{|S_{v_i}|}{|S|} \log_2 \frac{|S_{v_i}|}{|S|} = - \sum_{v_i \in Values(A)} P(v_i) \log_2 P(v_i)$$

$$GainRatio(S, A) = \frac{Gain(S, A)}{SplitInformation(S, A)}$$

26

## Gini Index

- ▶ Gini diversity index (used by CART)
  - ▶ Another measure that measures the impurity of a data set.
  - ▶  $S$  is a set of training examples associated with a node
  - ▶ Suppose there are  $n$  classes:  $C_i$  ( $i = 1, \dots, n$ )
  - ▶  $P(C_i)$  is the probability of examples in  $S$  that belong to  $C_i$ .
  - ▶ The Gini impurity of  $S$  with respect to classes can be measured as:

$$i(S) = \sum_{j \neq i} P(C_j) P(C_i) = 1 - \sum_{j=1}^n (P(C_j))^2$$

- ▶ Similar to entropy
  - ▶ Minimized if classes for all examples are the same
  - ▶ Maximized if equal proportion of classes
- ▶ Selection of attribute using Gini index selects an attribute  $A$  that most reduces the impurity due to partitioning on  $A$ :

$$\Delta i(S, A) = i(S) - \sum_{v_i \in Values(A)} \frac{|S_{v_i}|}{|S|} i(S_{v_i})$$

27

## Outline

- ▶ Overview of classification
- ▶ Basic concepts in decision tree learning
  - ▶ Data representation in decision tree learning
  - ▶ What is a decision tree?
    - ▶ Decision tree representation
- ▶ How to learn a decision tree from data
  - ▶ Basic decision tree learning algorithm
  - ▶ How to select best attribute
  - ▶ *Pruning decision tree*
- ▶ Other issues involved in decision tree learning

28

## Basic Decision Tree Learning Algorithm (Review)

1. Select the “best” attribute *A* for the root node
2. Create new descendents of the node according to the values of *A*:
3. Sort training examples to the descendent nodes.
4. For each descendent node,
  - ▶ if the training examples associated with the node belong to the same class, the node is marked as a leaf node and labeled with the class
  - ▶ else if there are no remaining attributes on which the examples can be further partitioned, the node is marked as a leaf node and labeled with the most common class among the training cases for classification;
  - ▶ else if there is no example for the node, the node is marked as a leaf node and labeled with the majority class in its parent node.
  - ▶ otherwise, recursively apply the process on the new node.

*when to  
terminate  
the  
recursive  
process*

29

## Overfitting Problem

- ▶ When to declare a node terminal in the basic algorithm:
  - ▶ grow each branch of the tree just deeply enough to classify the training examples as perfectly as possible.
- ▶ This strategy leads to producing deep branches that cover very few examples.
- ▶ This kind of trees *overfits* the training data in the following two situations:
  - ▶ there is noise in the data → tree fits the noise.
  - ▶ the number of training examples is too small to produce a representative sample of the true target function → tree is too specific to classify future examples well.

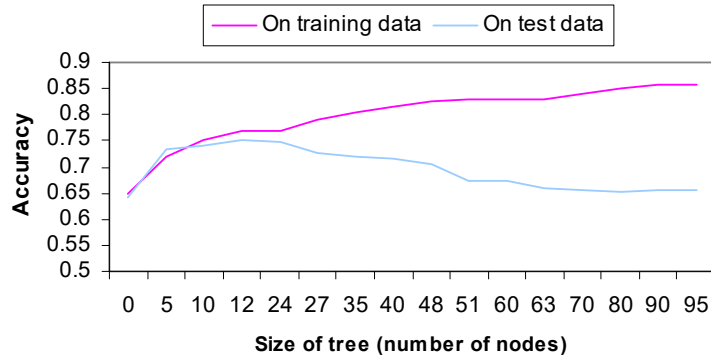
30

## Overfitting Problem (*Cont'd*)

- ▶ A definition of *overfitting*
  - ▶ Consider the error of a model (e.g., a tree)  $h$  over
    - ▶ training data:  $error_{train}(h)$
    - ▶ entire distribution  $D$  of data:  $error_D(h)$
  - ▶ Model  $h$  overfits training data if there is an alternative model  $h'$  such that
$$error_{train}(h) < error_{train}(h') \text{ and } error_D(h) > error_D(h')$$

31

## Overfitting in Decision Tree Learning (An Example Diagram)



Accuracy = 1 – error rate

32

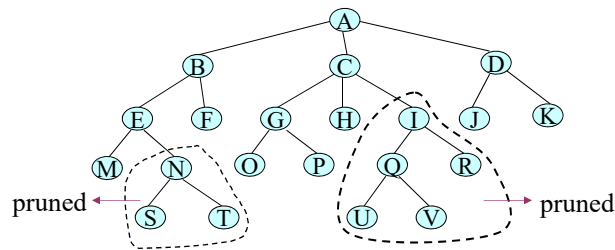
## Preventing Overfitting

- ▶ **Pre-pruning:** stop growing the tree when data split is not statistically significant
  - ▶ For example,
    - ▶ set a threshold  $\alpha > 0$  and declare a node terminal if percentage of examples in the most common class  $> \alpha$ , or
    - ▶ set a threshold  $\beta > 0$  and declare a node terminal if highest information gain  $< \beta$
  - ▶ Problems:
    - ▶ Hard to set the threshold value
    - ▶ The splitting is either stopped too soon or continued too far depending on the threshold
  - ▶ Using more complicated stopping rule does not help

33

## Preventing Overfitting (Cont'd)

- ▶ **Post-pruning:** grow full tree, allow it to overfit the data, and then remove some subtrees



- ▶ More successful in practice
- ▶ Criterion is needed to determine what to prune

34

## Post-pruning Functions

- ▶ Post-pruning functions are needed to determine which part(s) of the tree should be pruned
- ▶ Several post-pruning functions, including:
  - ▶ Reduced-error (Quinlan, 83)
  - ▶ Error-complexity (CART, 84)
  - ▶ Pessimistic error (Quinlan, 86)
  - ▶ Minimum description length (MDL) (SLIQ by Mehta *et al.* 1996)
  - ▶ Minimum error (Niblett & Bratko, 86)
  - ▶ Critical value (Mingers, 87)
- ▶ There is no single best pruning algorithm

35

## Reduced-Error Pruning

### ► Procedure

- Split training data into *growing* and *pruning/validation* sets
- Generate an overfitted decision tree using the *growing* set
- Post-pruning: Do until further pruning is harmful:
  - Consider each of the internal non-root nodes in the tree to be candidates for pruning
    - Prune a node by removing subtree rooted at this node, making it a leaf node, and assigning it the most common classification of the training examples affiliated with this node
    - Evaluate impact of pruning this node on the *pruning* set by
      - calculating the classification error rate of the pruned tree on the *pruning* set and comparing it with the error rate of the unpruned tree.
  - Greedily remove the one whose removal most reduces the error on *pruning* set.

### ► Aim to produce smallest version of most accurate subtree

### ► What if data is limited?

But may be sub-optimal. 36

## Error-Complexity Pruning

### ► Similar procedure to the reduced-error pruning

- Split training data into *growing* and *pruning/validation* sets
- Generate the decision tree with the *growing* set
- Post-pruning: Do until further pruning is harmful:
  - Evaluate impact of pruning each subtree on the *pruning* set
  - Greedily remove the one whose removal (resulting in tree  $T$ ) minimizes the following expression on the *pruning* set:

$$R_{\alpha}(T) = E(T) + \alpha L(T)$$

where  $E(T)$  is the classification error of tree  $T$ ,  $L(T)$  is the number of leaf nodes in  $T$  and  $\alpha$  is the complexity cost per leaf node.

- $R_{\alpha}(T)$  is a linear combination of the classification error rate of the tree and its complexity.

### ► Similar to reduced-error pruning, not suitable if the size of training data is small.

37

## Pessimistic-Error Pruning

- ▶ This method does not require a separate pruning/validation set.
- ▶ Suppose a subtree  $T_s$  contains  $L$  leaves and  $J$  of training examples associated with  $T_s$  are misclassified
- ▶ If we replace  $T_s$  with a leaf which misclassifies  $E$  of the associated training examples, the pruned tree will be accepted if

$$E + 1/2 < J + L/2$$

- ▶ Perform top-down traversing over internal nodes
  - ▶ If an internal node is pruned, all its descendants are removed
  - ▶ Relatively fast pruning
- ▶ Advantage
  - ▶ fast: no need to separate growing and pruning sets.
  - ▶ Tree building makes use of all the training data.

38

## Outline

- ▶ Overview of classification
- ▶ Basic concepts in decision tree learning
  - ▶ Data representation in decision tree learning
  - ▶ What is a decision tree?
    - ▶ Decision tree representation
- ▶ How to learn a decision tree from data
  - ▶ Basic decision tree learning algorithm
  - ▶ How to select best attribute
  - ▶ Pruning decision tree
- ▶ *Other issues involved in decision tree learning*

39



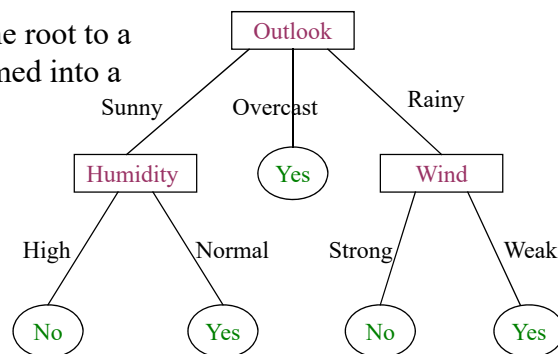
## Some Other Issues in Decision Tree Learning

- ▶ Convert decision trees to a set of rules
- ▶ How to deal with continuous attributes
- ▶ How to scale up decision tree learning
- ▶ Look-ahead approaches
  - ▶ Search for best sequence of individual tests.
- ▶ Multiple attributes per test
  - ▶ These approaches tend to have greatly increased complexities (i.e., larger search spaces).

40

## Convert Decision Trees to a Set of Rules

- ▶ Each branch from the root to a leaf can be transformed into a if-then rule.



- ▶ If (Outlook is Sunny and Humidity is High), then class is No.
- ▶ If (Outlook is Sunny and Humidity is Normal), then class is Yes.
- ▶ If (Outlook is Overcast), then class is Yes.
- ▶ .....

41

## How to deal with continuous attributes

Two ways:

- ▶ Discretization before learning decision tree
  - ▶ Converts a continuous attribute into a discrete attribute by partitioning the range of the continuous attribute into intervals.
  - ▶ Interval labels can then be used to replace actual data values.
- ▶ Dynamic discretization
  - ▶ Dynamically split the value range into two sub-ranges during the tree learning process

42

## Dynamic Discretization

- ▶ Dynamically split the value range into two sub-ranges and each descendent node corresponds to a sub-range.
  - ▶ Choose to split at the middle value between two examples which are in different categories

Temperature	PlayTennis?
40	No
48	No
60	Yes
70	Yes
80	Yes
90	No

The possible split points are  $\frac{48 + 60}{2} = 54$  and  $\frac{80 + 90}{2} = 85$

- ▶ Evaluate the binary splitting points using the splitting criterion used for selecting attribute. For example, choose the splitting point that leads to the best information gain.

43

## Dynamic Discretization

Example:

Temperature	PlayTennis?
40	No
48	No
60	Yes
70	Yes
80	Yes
90	No

The possible split points are  $\frac{48 + 60}{2} = 54$  and  $\frac{80 + 90}{2} = 85$

Possible splits:



Choose the one that gives the better information gain to compete with other attributes

44

## Scalable Decision Tree Learning Methods

- ▶ Scalability: deal with millions of examples and hundreds of attributes with reasonable speed
- ▶ Most algorithms assume data can fit in memory.
- ▶ Data mining research contributes to the scalability issue, especially for decision trees.
- ▶ Successful examples
  - ▶ SLIQ (Mehta *et al.*, 1996)
  - ▶ SPRINT (Shafer *et al.*, 1996)
  - ▶ RainForest (Gehrke, *et al.*, 1998)

45

## Some Other Issues in Decision Tree Learning

- ▶ Convert decision trees to a set of rules
- ▶ How to deal with continuous attributes
- ▶ How to scale up decision tree learning
- ▶ *Look-ahead approaches*
  - ▶ *Search for best sequence of individual tests.*
- ▶ *Multiple attributes per test/node*
  - ▶ *These approaches tend to have greatly increased complexities (i.e., larger search spaces).*

46

## Decision Trees: Strengths

- ▶ Comprehensibility: Small trees are highly interpretable and intuitive for humans
- ▶ Automatic attribute selection.
- ▶ Fast classification
- ▶ Relatively fast induction
- ▶ Applicable to both regression and classification problems.
- ▶ Mature technology

47

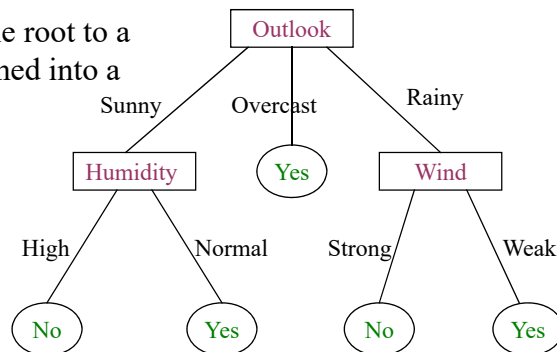
## Decision Trees: Weaknesses

- ▶ Trees can become incomprehensible when their size grows.
  - ▶ Rules converted from a tree
    - ▶ are *mutually exclusive*,
    - ▶ share at least one attribute (the root)
  - ▶ Thus, the size of the tree can grow much larger than the logic needed for overlapping rules.
  - ▶ As an example, in a successful application of ID3 to a chess end game (Mitchie, 86), the tree representation could not be understood at all even by the chess experts.

48

## Convert Decision Trees to a Set of Rules

- ▶ Each branch from the root to a leaf can be transformed into a if-then rule.



- ▶ If (Outlook is Sunny and Humidity is High), then class is No.
- ▶ If (Outlook is Sunny and Humidity is Normal), then class is Yes.
- ▶ If (Outlook is Overcast), then class is Yes.
- ▶ .....

49

## Decision Trees: Weaknesses (*Cont'd*)

- ▶ When using only one attribute at each internal node, trees are limited to axis-parallel partitions of the instance space
- ▶ Instability
  - ▶ If we change the data a little, the tree can change a lot. So does accuracy.
  - ▶ Can be improved with: Random forests

50

## Random Forests

- ▶ Grow a forest of many trees (e.g., 500)
- ▶ Grow each tree on an independent *bootstrap sample* from the training data.
  - ▶ Bootstrap sample: Sample  $N$  examples at random with replacement
- ▶ At each node:
  - ▶ Select  *$m$  attributes at random* out of all  $M$  possible attributes (independently for each node).
  - ▶ Find the best split on the selected  *$m$*  variables.
- ▶ Grow the trees to maximum depth
- ▶ Vote/average the trees to get predictions for new data

51

## Random Forests

- ▶ Advantages:

- ▶ Accuracy

- ▶ Random Forests is competitive with the best known machine learning methods

- ▶ Stability in predictive performance

- ▶ If we change the data a little, the individual trees may change but the forest is relatively stable because it is a combination of many trees.

- ▶ Disadvantage

- ▶ Interpretability decreases – not interpretable if there are, say, 500 trees.

52

## Summary of Decision Tree Learning

- ▶ Decision tree represents classification knowledge

- ▶ Decision tree learning is a top-down recursive partitioning process. It is a two phase process if post-pruning is used:

- ▶ Tree building

- ▶ An attribute selection criterion (also called splitting criterion) is used: information gain, gain ratio, gini index, etc.

- ▶ Post-pruning tree

- ▶ Reduced-error (Quinlan, 87)
    - ▶ Error-complexity (CART, 84)
    - ▶ Pessimistic error (Quinlan, 86)

- ▶ Some other issues

53

## Next Class

- ▶ Data Preprocessing (Chapter 3)