

Decision tree classifier

Prof. E.P.Ephzibah

General Approach to Classification

- ▶ The data classification process:
- ▶ (a) Learning: Training data are analyzed by a classification algorithm. The learned model or classifier is represented in the form of classification rules (if -then rules).
- ▶ (b) Classification: Test data are used to estimate the accuracy of the classification rules. If the accuracy is considered acceptable, the rules can be applied to the classification of new data tuples.

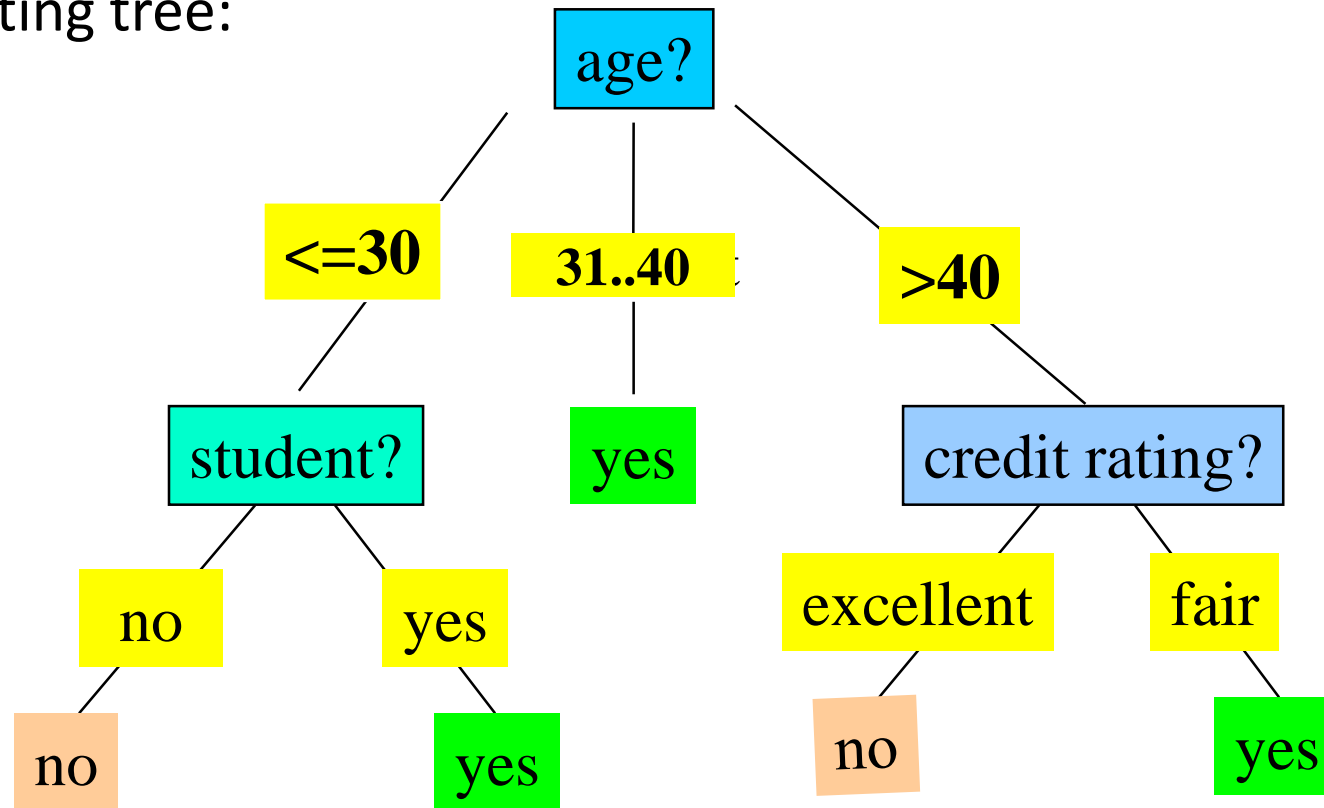
Decision Tree Induction: An Example

□ Training data set: 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

Decision Tree Induction: An Example

□ Resulting tree:



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)

- ID3 uses information gain as its attribute selection measure.
- 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 $|C_{i,D}|/|D|$ $Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$
- Expected information (entropy) needed to classify a tuple in D :
$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$
- Information needed (after using A to split D into v partitions) to classify D :

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

- The information gained by branching on attribute A

Attribute Selection: Information Gain

- Class P: buys_computer = “yes”
- Class N: buys_computer = “no”

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

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

Attribute Selection: Information Gain

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

$$\begin{aligned} \text{Info}_{\text{age}}(D) &= \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) \\ &\quad + \frac{5}{14} I(3,2) = 0.694 \end{aligned}$$

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

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.246$$

Similarly,

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

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

$$\text{Gain}(\text{credit_rating}) = 0.048$$

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 is selected as the split-point for A
- ▶ Split:
 - ▶ D1 is the set of tuples in D satisfying $A \leq \text{split-point}$, and D2 is the set of tuples in D satisfying $A > \text{split-point}$

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)

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

- ▶ $GainRatio(A) = Gain(A)/SplitInfo(A)$
- ▶ Ex.
$$SplitInfo_{income}(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \times \log_2\left(\frac{4}{14}\right) = 1.557$$
- ▶ $gain_ratio(income) = 0.029/1.557 = 0.019$
- ▶ The attribute with the maximum gain ratio is selected as the splitting attribute

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

- ▶ 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:

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

- ▶ The attribute provides the smallest $gini_{split}(D)$ (or the largest reduction in impurity) is chosen to split the node
(*need to enumerate all the possible splitting points for each attribute*)

Computation of Gini Index

- ▶ Ex. D has 9 tuples in buys_computer = “yes” and 5 in “no”

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

- ▶ Suppose the attribute income partitions D into 10 in D_1 : {low, medium} and 4 in D_2

$$\begin{aligned} gini_{income \in \{low, medium\}}(D) &= \left(\frac{10}{14}\right) Gini(D_1) + \left(\frac{4}{14}\right) Gini(D_2) \\ &= \frac{10}{14} \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right) \\ &= 0.443 \\ &= Gini_{income \in \{high\}}(D). \end{aligned}$$

$Gini_{\{low, high\}}$ is 0.458; $Gini_{\{medium, high\}}$ is 0.450. Thus, split on the {low, medium} (and {high}) since it has the lowest Gini index

- ▶ All attributes are assumed continuous-valued
- ▶ May need other tools, e.g., clustering, to get the possible split values
- ▶ Can be modified for categorical attributes

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

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

- Build a **decision tree** for the given training data in the table (Buy Computer data), predict the class of the following new

example :

- age<=30, income=medium, student=yes, credit-rating=fair**