



Decision Tree Induction

Outline

- ❑ Decision Tree Induction: Basic Idea and Algorithm
- ❑ Alternative Attribute Selection Measures in Decision Tree Induction
- ❑ Overfitting and Tree Pruning
- ❑ Decision Tree Construction in Large Datasets
- ❑ Visualization of Decision Trees and Tree Construction by Visual Data Mining



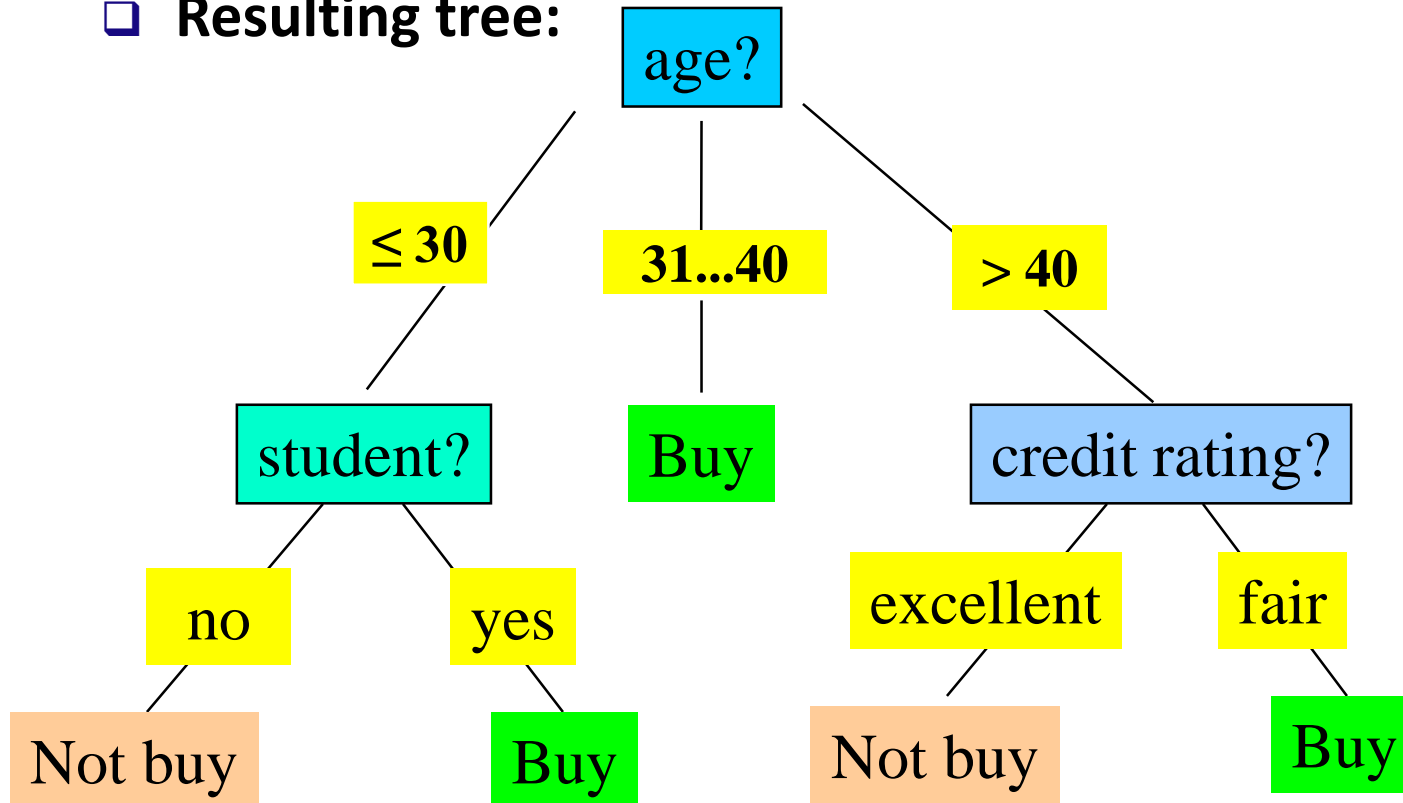
Decision Tree Induction: Basic Idea and Algorithm

Decision Tree Induction: An Example

Decision tree construction:

- A top-down, recursive, divide-and-conquer process

Resulting tree:



Training data set: Who 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 |

Note: The data set is adapted from
“Playing Tennis” example of R. Quinlan

From Entropy to Info Gain: A 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 $\{y_1, y_2, \dots, y_m\}$

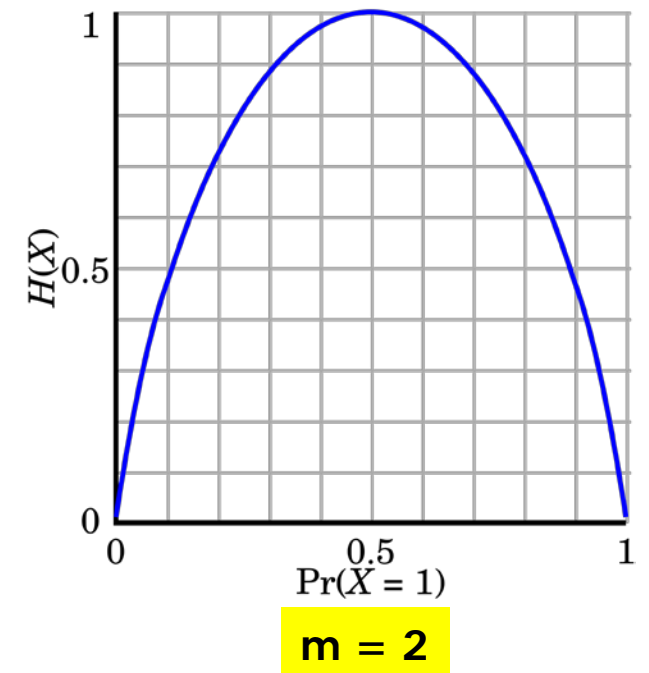
$$H(Y) = - \sum_{i=1}^m p_i \log(p_i) \quad \text{where } p_i = P(Y = y_i)$$

□ Interpretation

- Higher entropy \rightarrow higher uncertainty
- Lower entropy \rightarrow lower uncertainty

□ Conditional entropy

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



Information Gain: An Attribute Selection Measure

- ❑ Select the attribute with the highest information gain (used in typical decision tree induction algorithm: ID3)
- ❑ 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 :

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

Example: Attribute Selection with 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 | p _i | n _i | I(p _i , n _i) |
|---------|----------------|----------------|-------------------------------------|
| <=30 | 2 | 3 | 0.971 |
| 31...40 | 4 | 0 | 0 |
| >40 | 3 | 2 | 0.971 |

| 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 |

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

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

Hence

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

Similarly, we can get

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$

Decision Tree Induction: Algorithm

□ Basic algorithm

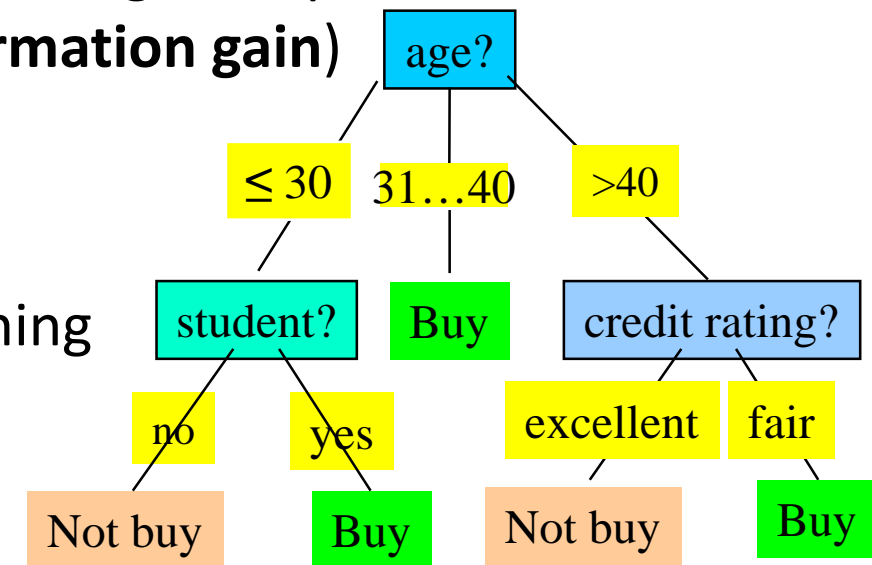
- Tree is constructed in a **top-down, recursive, divide-and-conquer manner**
- At start, all the training examples are at the root
- Examples are partitioned recursively based on selected attributes
- On each node, attributes are selected based on the training examples on that node, and a heuristic 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
- There are no samples left

□ Prediction

- **Majority voting** is employed for classifying the leaf



How to Handle Continuous-Valued Attributes?

- ❑ Method 1: Discretize continuous values and treat them as categorical values

 - ❑ E.g., age: < 20, 20...30, 30...40, 40...50, > 50

- ❑ Method 2: Determine the **best split point** for continuous-valued attribute A

 - ❑ Sort the value A in increasing order, E.g., 15, 18, 21, 22, 24, 25, 29, 31, ...

 - ❑ *Possible split point*: The midpoint between *each pair of adjacent values*

 - ❑ $(a_i + a_{i+1})/2$ is the midpoint between the values of a_i and a_{i+1}

 - ❑ e.g., $(15+18)/2 = 16.5, 19.5, 21.5, 23, 24.5, 27, 30, \dots$

 - ❑ The point with the *maximum information gain* for A is selected as the **split-point** for A

- ❑ Split: Based on split point P

 - ❑ The set of tuples in D satisfying $A \leq P$ vs. those with $A > P$