# CS 412 Intro. to Data Mining

## Chapter 8. Classification: Basic Concepts

Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017

# Example: Attribute Selection with Information Gain

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

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

| 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'es and 3 no's.

Hence

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

Similarly, we can get

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit\_rating) = 0.048$$

12

# Decision Tree Induction: Algorithm

- ❑ Basic algorithm
  - ❑ Tree is constructed in a **top-down, recursive, divide-and-conquer** manner

    *แบ่ง data จนกว่าจะ...*
  - ❑ At start, all the training examples are at the root
  - ❑ Examples are partitioned recursively based on selected attributes

    *data แบวพร้อมเดิมไปเรื่อยๆ*
  - ❑ 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

13

# 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, …
  - The point with the *maximum information gain* for A is selected as the **split-point** for A
- Split: Based on split point P
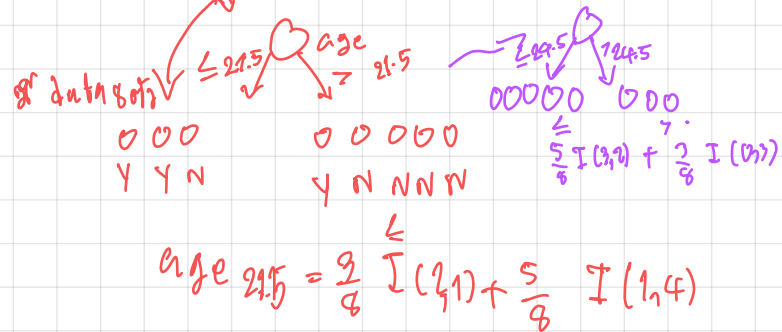  - The set of tuples in D satisfying A ≤ P vs. those with A > P

# Math 1   Categorize

15, 18, 21, 22, 24, 25, 29, 31, ...

$< 16, 16-22, 22-30, > 31$

## math 2   Best splitpoint

Best data    15, 18, 21, 22, 24, 25, 29, 31, ...



gr data sets   $\leq 21.5$  age  $> 21.5$

Y Y N          Y N N N N

$\leq 29.5$   $\geq 24.5$

OOOOO   OOO

$\frac{5}{8} I(3,2) + \frac{3}{8} I(3,0)$

$age_{21.5} = \frac{3}{8} I(2,1) + \frac{5}{8} I(1,4)$

# Gain Ratio: A Refined Measure for Attribute Selection

- ❑ Information gain measure is biased towards attributes with a large number of values

- ❑ Gain ratio: Overcomes the problem (as a normalization to information gain)

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

- ❑ GainRatio(A) = Gain(A)/SplitInfo(A)

- ❑ The attribute with the maximum gain ratio is selected as the splitting attribute

- ❑ Gain ratio is used in a popular algorithm C4.5 (a successor of ID3) by R. Quinlan

- ❑ Example

- ❑ $SplitInfo_{income}(D) = -\frac{4}{14}\log_2\frac{4}{14} - \frac{6}{14}\log_2\frac{6}{14} - \frac{4}{14}\log_2\frac{4}{14} = 1.557$

- ❑ GainRatio(income) = 0.029/1.557 = 0.019

# Another Measure: Gini Index

❑ Gini index: Used in CART, and also in IBM IntelligentMiner

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

  ❑ $gini(D) = 1 - \sum_{j=1}^{n} p_j^2$    $\geq p \log p \to (-p \log p) + (-p \log p)$

    ❑ $p_j$ 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*)

*การตัดแต่งกิ่ง*
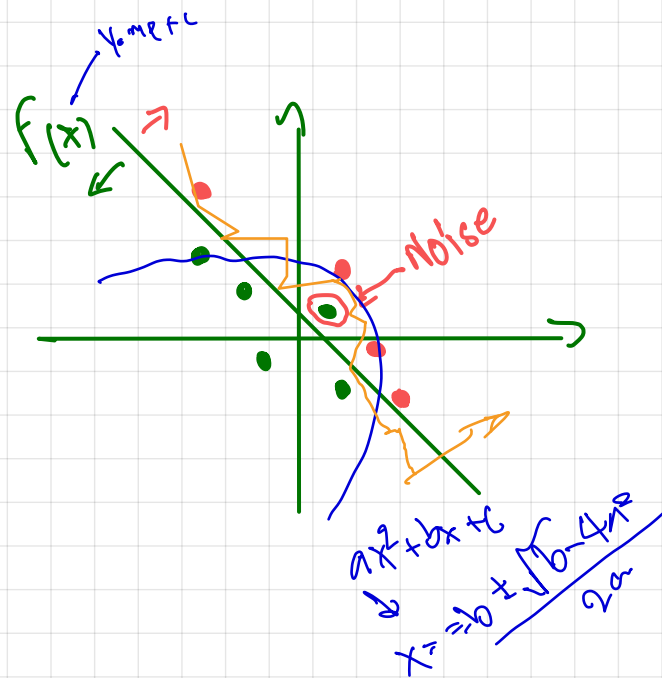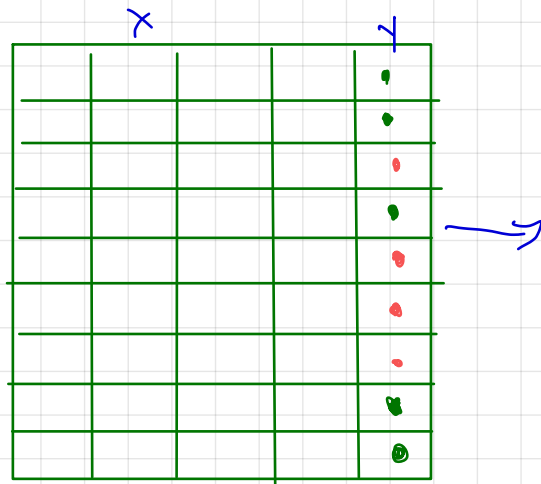
❑ <u>Overfitting</u>: An induced tree may overfit the training data

 ❑ Too many branches, some may reflect anomalies due to noise or outliers

 ❑ Poor accuracy for unseen samples

❑ Two approaches to avoid overfitting

 ❑ <u>Prepruning</u>: *Halt tree construction early*-do not split a node if this would result in the goodness measure falling below a threshold

*หยุดการโตของต้นไม้ ก่อน*

  ❑ Difficult to choose an appropriate threshold

 ❑ <u>Postpruning</u>: *Remove branches* from a "fully grown" tree—get a sequence of progressively pruned trees

*ตัดกิ่งออกหลังไปโตจนสุด แล้ว*

  ❑ Use a set of data different from the training data to decide which is the "best pruned tree"

20

X     Y

$f(x)$    $y_0 = mx + c$

Noise

$ax^2 + bx + c$

$x = \dfrac{-b \pm \sqrt{b^2 - 4ac}}{2a}$

Occam's Rezer