



# **CS 412 Intro. to Data Mining**

## **Chapter 8. Classification: Basic Concepts**

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





# **Chapter 8. Classification: Basic Concepts**

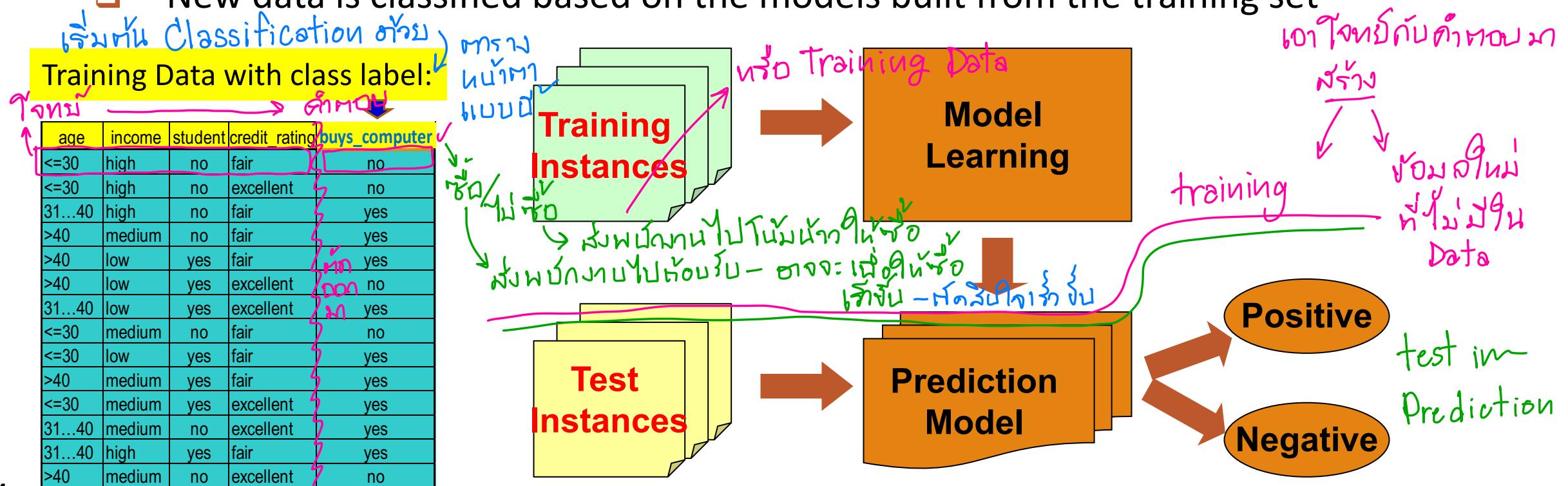
---

- Classification: Basic Concepts
- Decision Tree Induction
- Bayes Classification Methods
- Linear Classifier
- Model Evaluation and Selection
- Techniques to Improve Classification Accuracy: Ensemble Methods
- Additional Concepts on Classification
- Summary



# **Supervised vs. Unsupervised Learning (1)**

- **Supervised learning (classification)**
    - Supervision: The training data such as observations or measurements are accompanied by **labels** indicating the classes which they belong to
    - New data is classified based on the models built from the training set

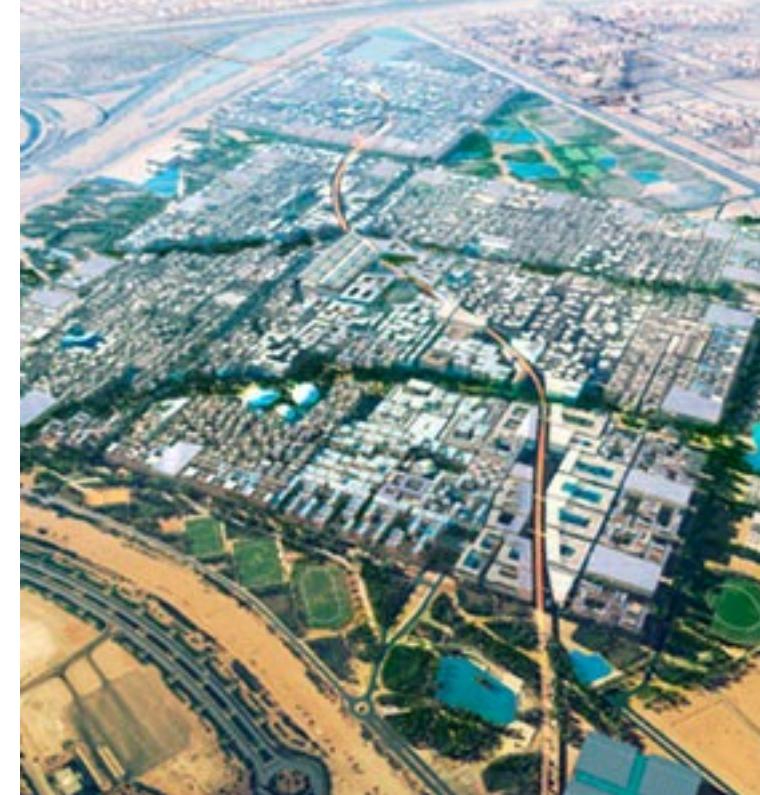
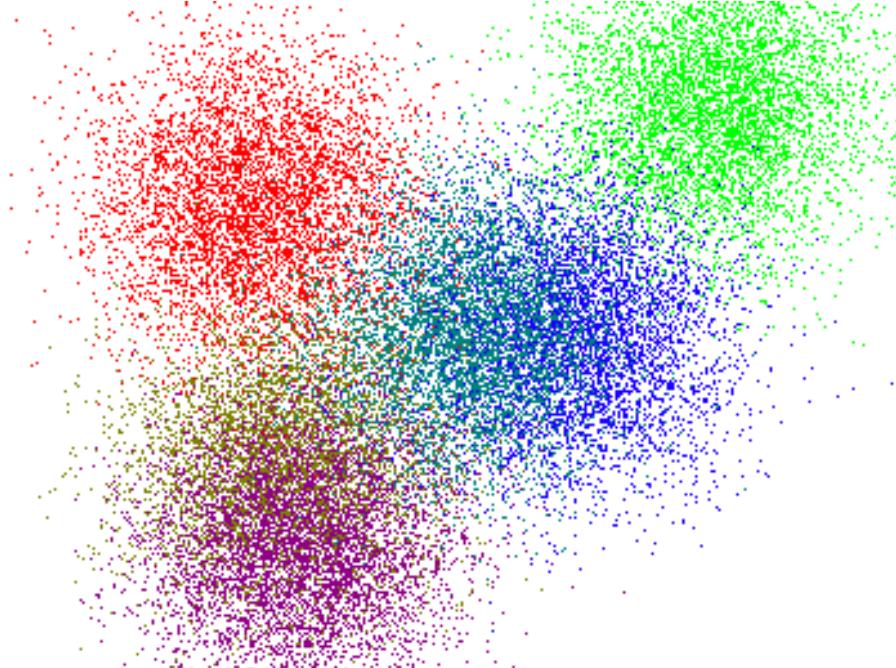


$x$ : feature မျက်နှာ  
 $y$ : target  
 $y$ : ဘဏ် /  $x$ : သွေထူး (training) → ပြည့်စုစုပါရို့  
 $x$ : မျက်နှာ  
 $y$ : မျက်နှာ

## Supervised vs. Unsupervised Learning (2)

Classification → Learning → ဂျာတဲ့ ဘဏ်ရေးလုပ်ငန်း

- Unsupervised learning (clustering) → မျက်နှာ ပြည့်စုစုပါရို့  
မျက်နှာ  $X$  ပြည့်စုစုပါရို့  $y$  } ဂေါ်ဆိုချုပ်လုပ်ငန်း
- The class labels of training data are unknown
- Given a set of observations or measurements, establish the possible existence of classes or clusters in the data

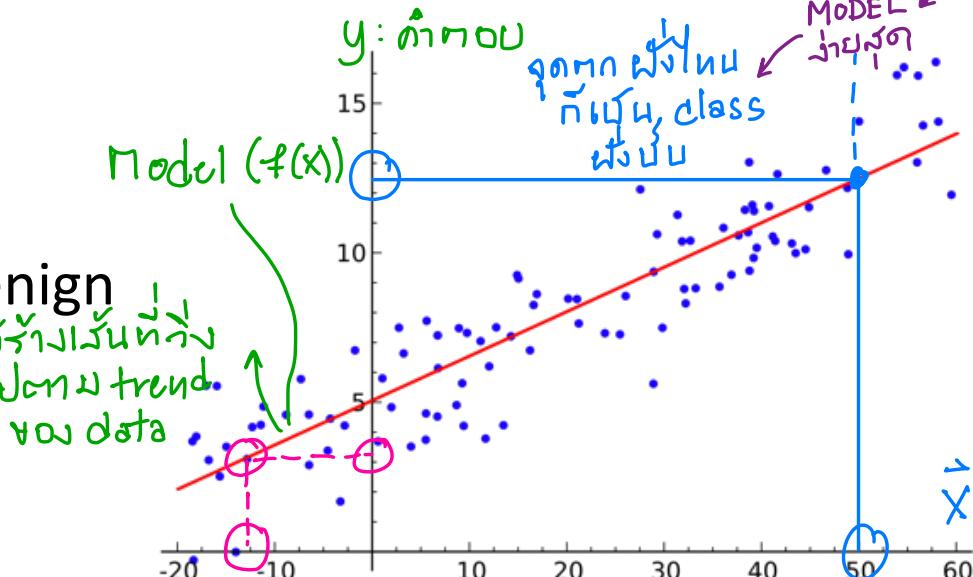
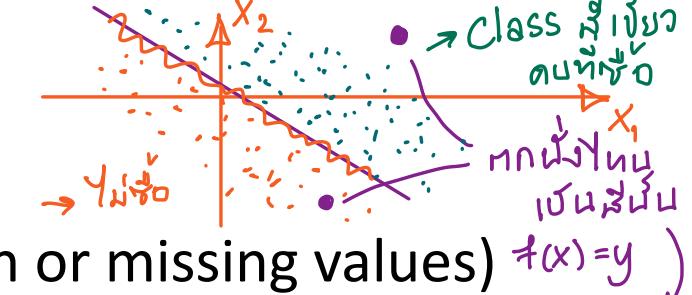


# Prediction Problems: Classification vs. Numeric Prediction

- Classification
  - ↳ ក្នុងព័ត៌មាន  
Ex. ទីតាំង/រូបថត → នគ. ហេតុបារាំង ធម៌=ប៊ូល
- Predict categorical class labels (discrete or nominal)
- Construct a model based on the training set and the **class labels** (the values in a classifying attribute) and use it in classifying new data
- Numeric prediction
- Model continuous-valued functions (i.e., predict unknown or missing values)
- Typical applications of classification
- Credit/loan approval
- Medical diagnosis: if a tumor is cancerous or benign
- Fraud detection: if a transaction is fraudulent
- Web page categorization: which category it is

Regression

$x_1, x_2, \dots$  = Feature  
ឈាមគោរព



# Classification—Model Construction, Validation and Testing

- **Model construction**
  - Each sample is assumed to belong to a predefined class (shown by the **class label**)
  - The set of samples used for model construction is **training set**
  - Model: Represented as decision trees, rules, mathematical formulas, or other forms
- **Model Validation and Testing:**
  - **Test:** Estimate accuracy of the model
    - The known label of test sample is compared with the classified result from the model
    - **Accuracy:** % of test set samples that are correctly classified by the model
    - Test set is independent of training set
  - **Validation:** If *the test set* is used to select or refine models, it is called **validation (or development) (test) set**
- **Model Deployment:** If the accuracy is acceptable, use the model to classify new data

# **Chapter 8. Classification: Basic Concepts**

---

- Classification: Basic Concepts
- Decision Tree Induction
- Bayes Classification Methods
- Linear Classifier
- Model Evaluation and Selection
- Techniques to Improve Classification Accuracy: Ensemble Methods
- Additional Concepts on Classification
- Summary



# Decision Tree Induction: An Example

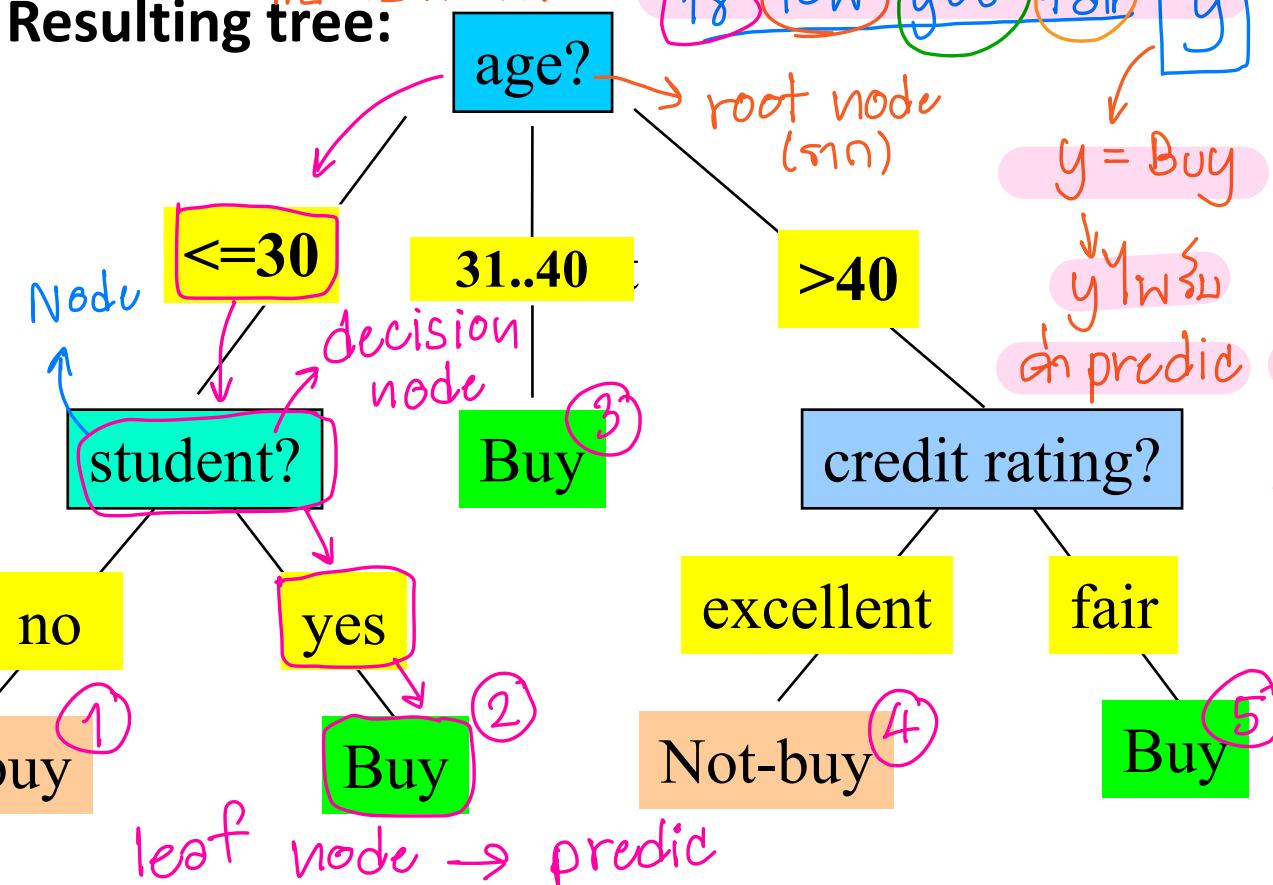
## □ Decision tree construction:

បន្ថែករាយការសេវា

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

ការសេវា

## □ Resulting tree:



Training data set: Who buys computer?

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

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

leaf Node

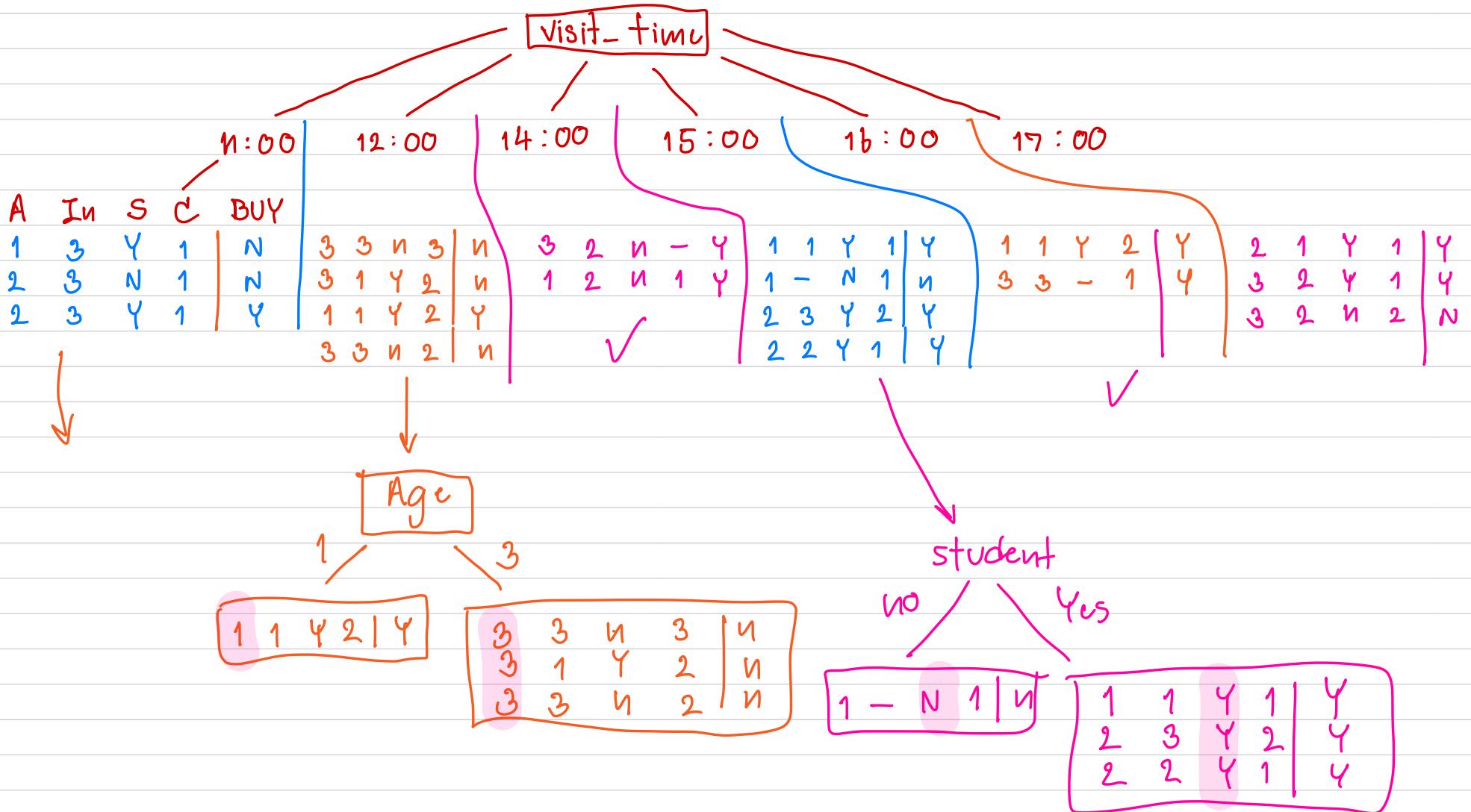
35 high No excellent  $\rightarrow y \rightarrow y' = \text{Buy} \rightarrow \text{Buy } ③$

30 median yes fair  $\rightarrow \text{Buy} \rightarrow ②$

15 low No excellent  $\rightarrow \text{Not-Buy} \rightarrow ①$

↗ តើ root node  $\rightarrow$  Data កំណែរបស់វា តើម្ខាត ពីរក្នុង

1	$L=30$	low	no	fair	no
2	$31 \dots 40$	medium	yes	ex	yes
3	$> 40$	high			



$f_1, f_2$

Target

A	A		Y
A	A		N
B	A		Y
A	B		N
B	B		N

$f_1$

$f_2$

A      B

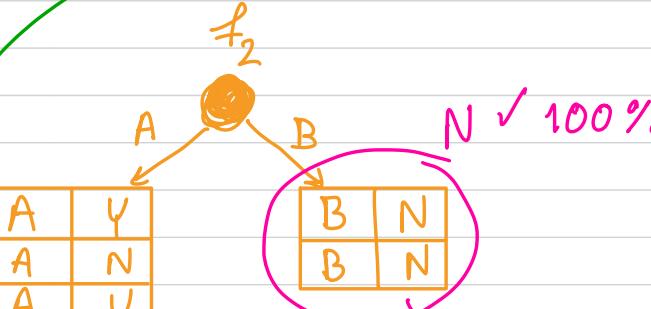
A	Y
A	N
A	N

B	Y
B	N

Initial Data:

\* \* \* A & common sense Qumsg

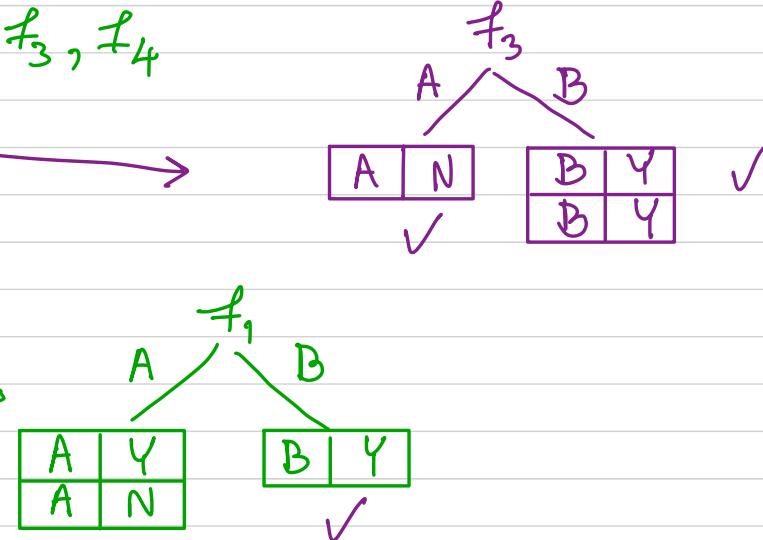


$f_1, f_3, f_4$

recursive

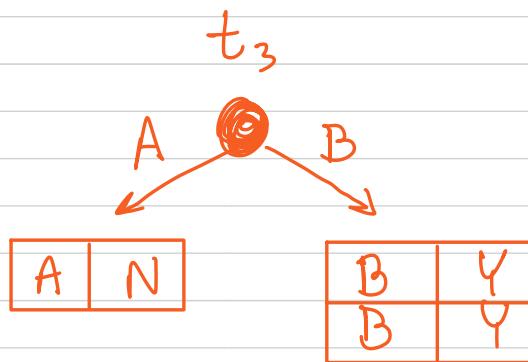
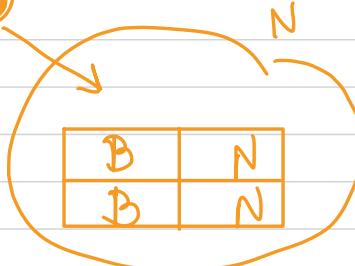
$f_1, f_3$

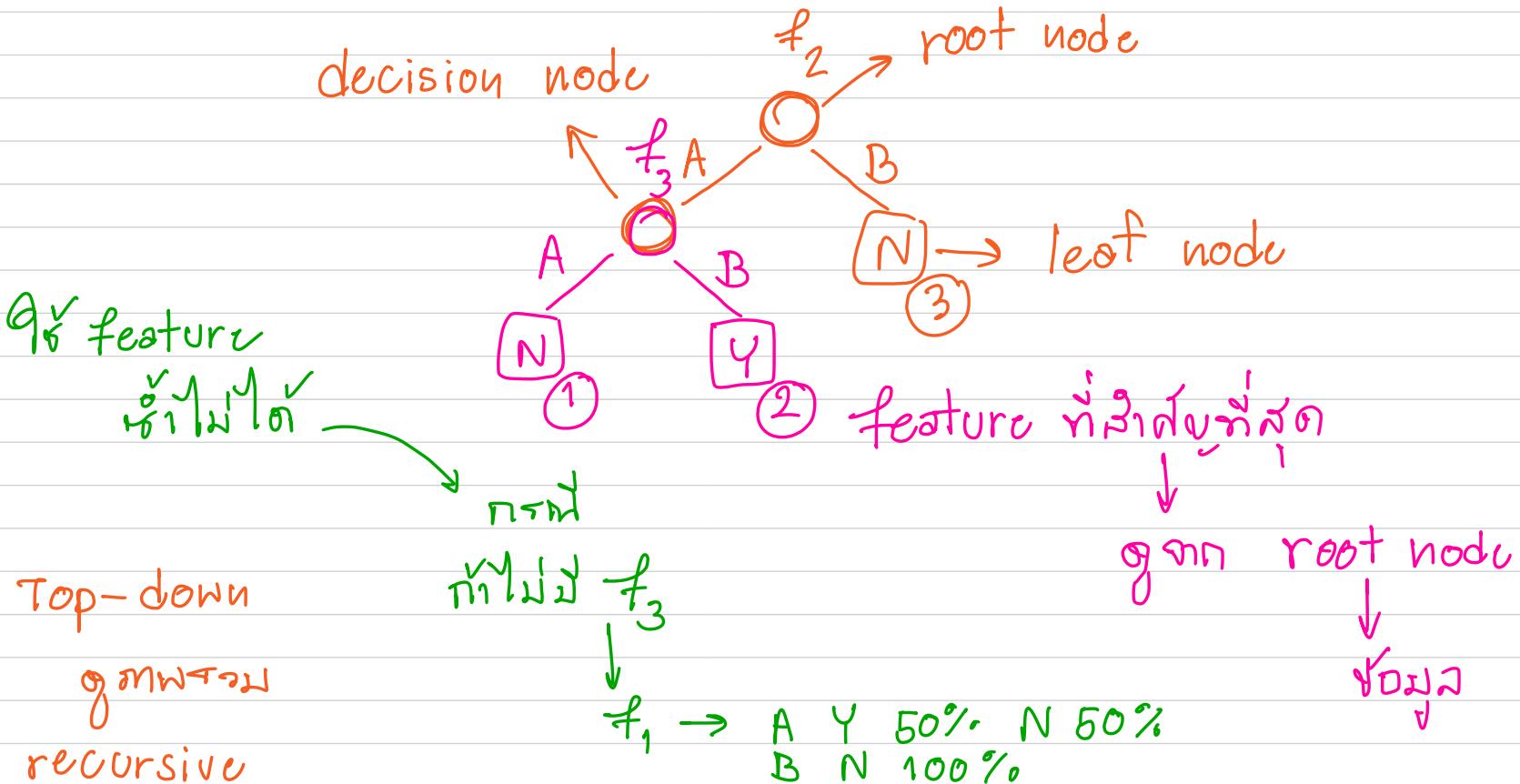
A	A	B		Y
A	A	A		N
B	A	B		Y



↓

$t_1$	$t_3$	$t_2$
A   A   B   Y		
A   A   A   N		
B   A   B   Y		





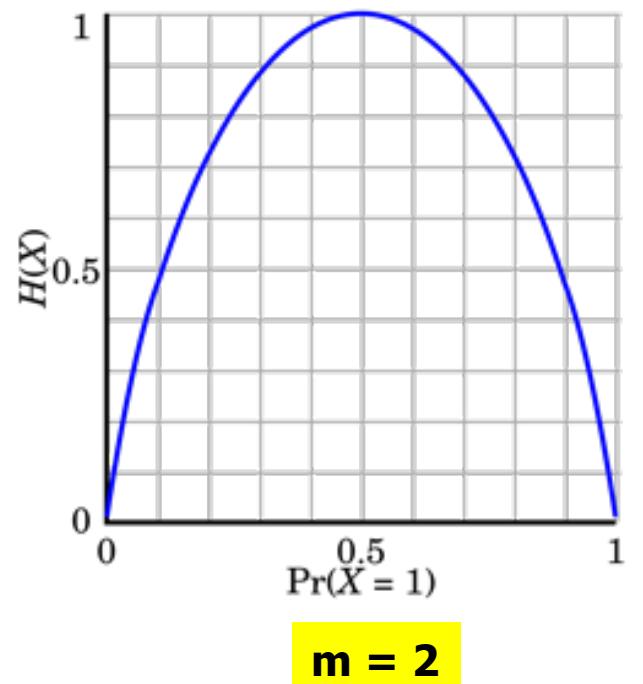
# From Entropy to Info Gain: A Brief Review of Entropy

- Entropy (Information Theory)
  - A measure of uncertainty associated with a random number
  - Calculation: For a discrete random variable  $Y$  taking  $m$  distinct values  $\{y_1, y_2, \dots, y_m\}$

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

- 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/C4.5)
  - 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:  
*class yes in sum នៃមុន្ត ជាដែល*

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

- 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

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

Node หัด  
↑ นี้

↑ ตัวอย่าง → เสือกเป็น root

มากที่สุด

↑ กฎแบ่งโดย feature A แล้ว

# Example: Attribute Selection with Information Gain

- Class P: buys\_computer = "yes" = 9
- Class N: buys\_computer = "no" = 5

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

age	$p_i$	$n_i$	$I(p_i, n_i)$
$\leq 30$	2	3	0.971
31...40	4	0	0
$>40$	3	2	0.971

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 ✓

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

feature

$$+ \frac{5}{14} I(3,2) = 0.694$$

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

Hence

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

Similarly, we can get

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

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

entropy  $\rightarrow$  ការមុនប្រាស់  
information  $\rightarrow$  ព័ត៌មាន

ឯកសារ

ក្រុងការបង្កើត  
root node

# 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	$p_i$	$n_i$	$I(p_i, n_i)$
$\leq 30$	2	3	0.971
$31 \dots 40$	4	0	0
$>40$	3	2	0.971

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

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

$\frac{5}{14} I(2,3)$  means “age  $\leq 30$ ” 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

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

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

$$\text{Info}(D) = I(2,3) = -\frac{3}{5} \log_2 \frac{2}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$\text{Info}_{\text{credit}}(D) =$$

$$\text{Info}_{\text{student}}(D) =$$

$$\text{Info}_{\text{credit}}(D) = \frac{3}{5} I(1,2) + \frac{2}{5} (1,1)$$

$$\begin{aligned} & \text{high } M \\ & \frac{2}{5} I(0,2) + \frac{3}{5} I(1,1) + \frac{2}{5} I(1,0) \end{aligned}$$

medium  
low

$\boxed{(2,0)}$   $\rightarrow$  1111000

Data for min Age  $\leq 30$

Student

no

yes

Age	Income	Credit_rating	buys_computer
<=30	high	fair	no
<=30	high	excellent	no
>30	high	fair	yes
<=30	low	fair	yes
<=30	low	fair	no
>30	low	excellent	yes
>30	low	excellent	no

Age	Income	Credit_rating	buys_computer
<=30	high	fair	yes
>30	low	excellent	no
>30	high	excellent	yes
<=30	low	fair	yes
<=30	low	fair	yes
<=30	low	excellent	yes
>30	high	fair	yes

Student : Yes

Age	Income	Credit_rating	buys_computer
<=30	high	fair	no
<=30	high	excellent	no
>30	high	fair	yes
<=30	low	fair	yes
<=30	low	fair	no
>30	low	excellent	yes
>30	low	excellent	no

Student : No

Age	Income	Credit_rating	buys_computer
<=30	high	fair	yes
>30	low	excellent	no
>30	high	excellent	yes
<=30	low	fair	yes
<=30	low	fair	yes
<=30	low	excellent	yes
>30	high	fair	yes

प्र० Student: no

$$\text{Info}_{\text{student: no}}(D) = I(3, 4) = -\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right) = 0.985$$

$$\begin{aligned} \text{Info}_{\text{age}}(D) &= \frac{4}{7}I(1, 3) + \frac{3}{7}I(2, 1) = \frac{4}{7}(0.811) + \frac{3}{7}(0.918) \\ &= 0.857 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{income}}(D) &= \frac{3}{7}I(1, 2) + \frac{4}{7}I(2, 2) = \frac{3}{7}(0.918) + \frac{4}{7}(1) \\ &= 0.965 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{credit-rating}}(D) &= \frac{3}{7}I(1, 2) + \frac{4}{7}I(2, 2) = \frac{3}{7}(0.918) + \frac{4}{7}(1) \\ &= 0.965 \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{age}) &= \text{Info}_{\text{student: no}} - \text{Info}_{\text{age}}(D) \\ &= 0.985 - 0.857 = 0.128 \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{income}) &= \text{Info}_{\text{student: no}} - \text{Info}_{\text{income}}(D) \\ &= 0.985 - 0.965 = 0.02 \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{credit\_rating}) &= \text{Info}_{\text{student: no}} - \text{Info}_{\text{credit\_rating}}(D) \\ &= 0.985 - 0.965 = 0.02 \quad \times \end{aligned}$$

age	$p_i$	$n_i$	$p_i+n_i$	$I(p_i+n_i)$
$\leq 30$	1	3	4	
$> 30$	2	1	3	

$$I(1, 3) = -\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right) = 0.811$$

$$I(2, 1) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = 0.918$$

income	$p_i$	$n_i$	$p_i+n_i$	$I(p_i+n_i)$
high	1	2	3	
low	2	2	4	

$$I(1, 2) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) = 0.918$$

$$I(2, 2) = -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) = 1$$

credit_rating	$p_i$	$n_i$	$p_i+n_i$	$I(p_i+n_i)$
excellent	1	2	3	
fair	2	2	4	

$$I(1, 2) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) = 0.918$$

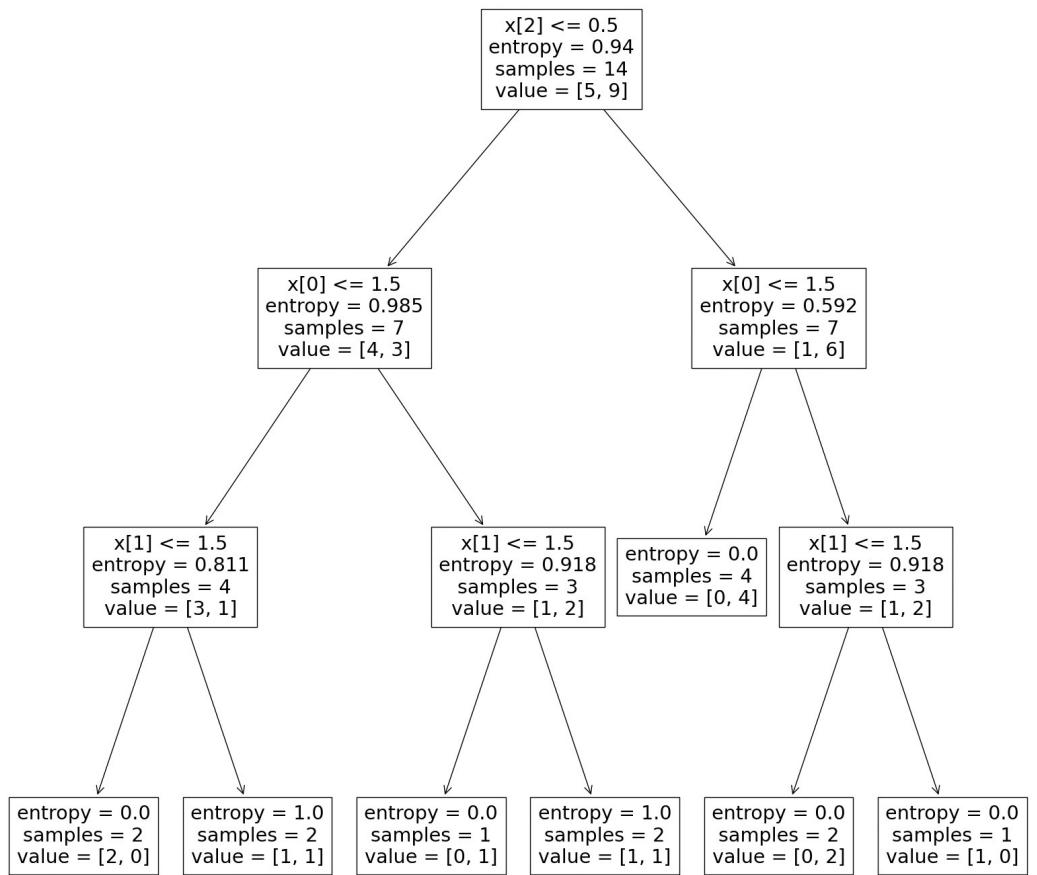
$$I(2, 2) = -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) = 1$$

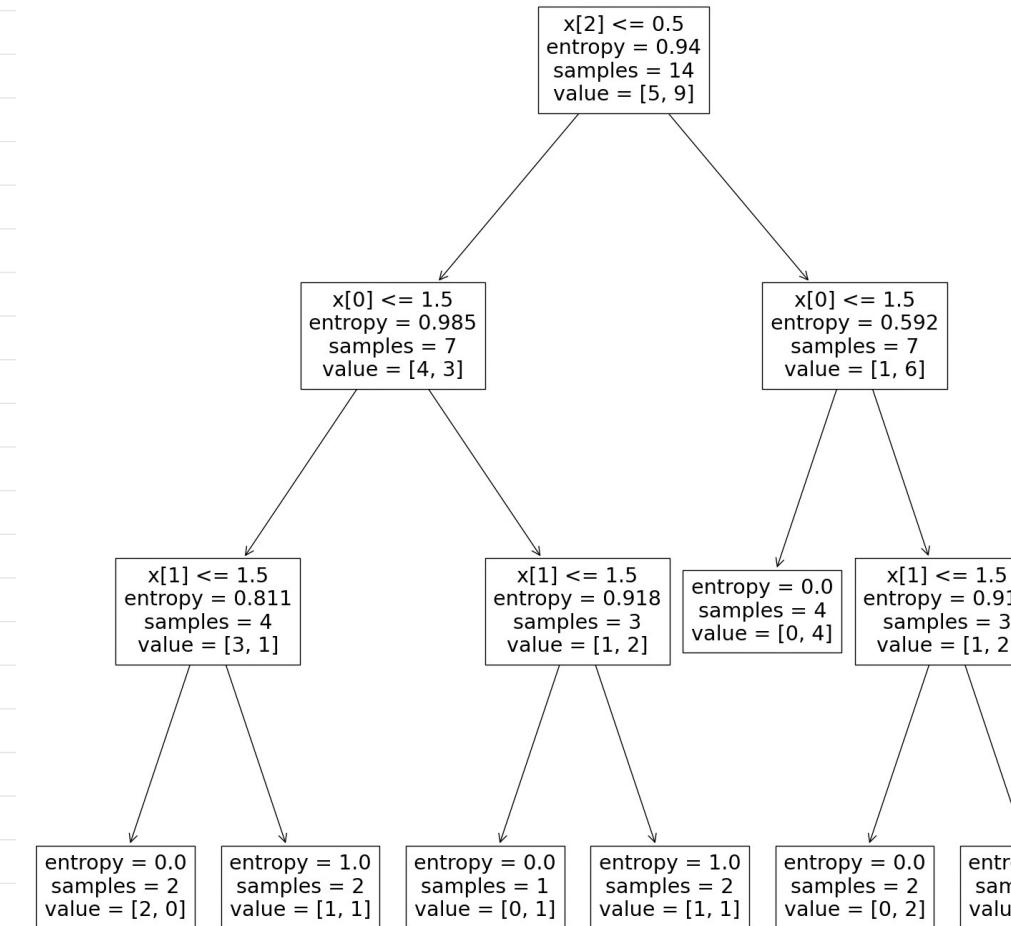
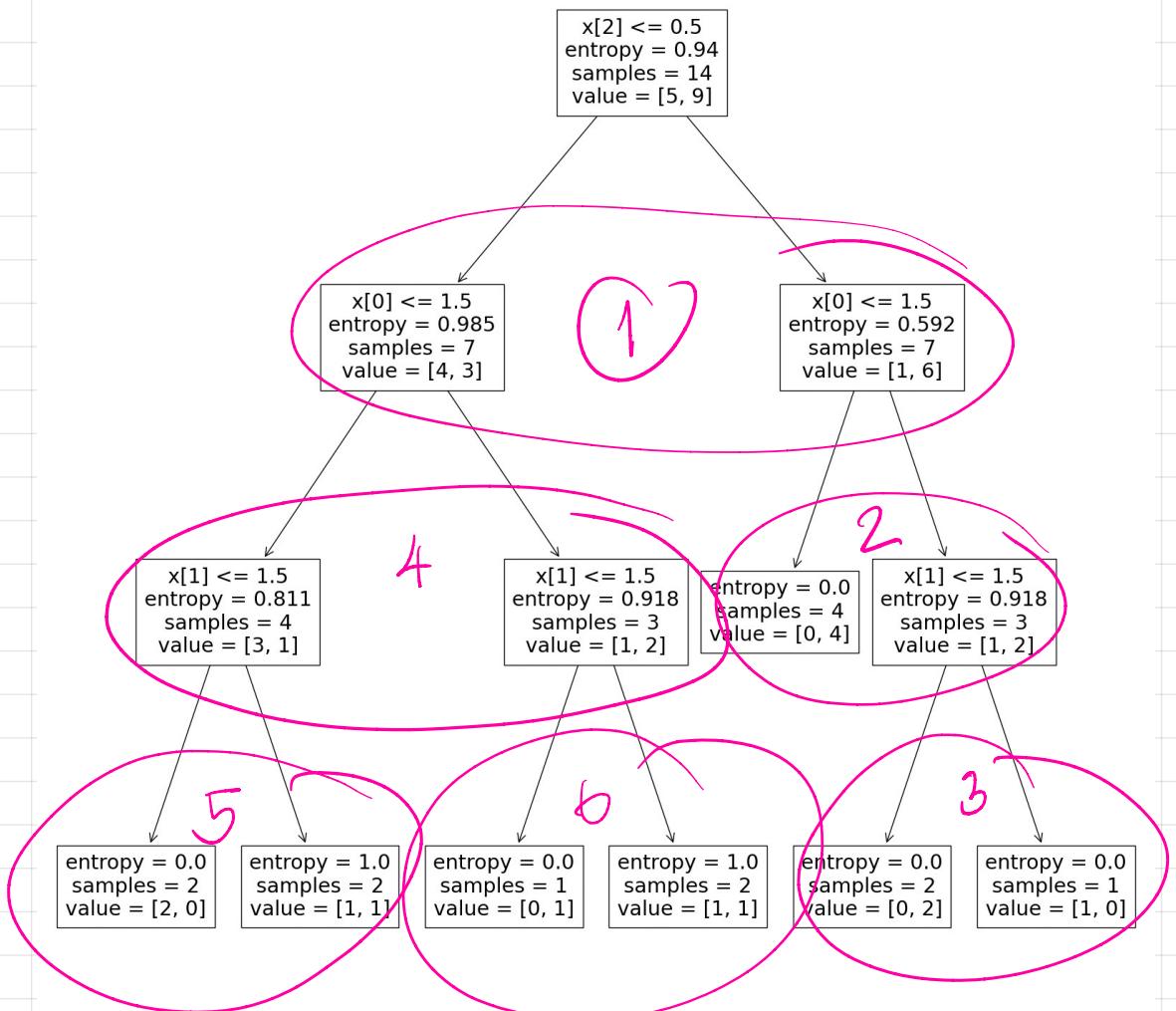
Student : No

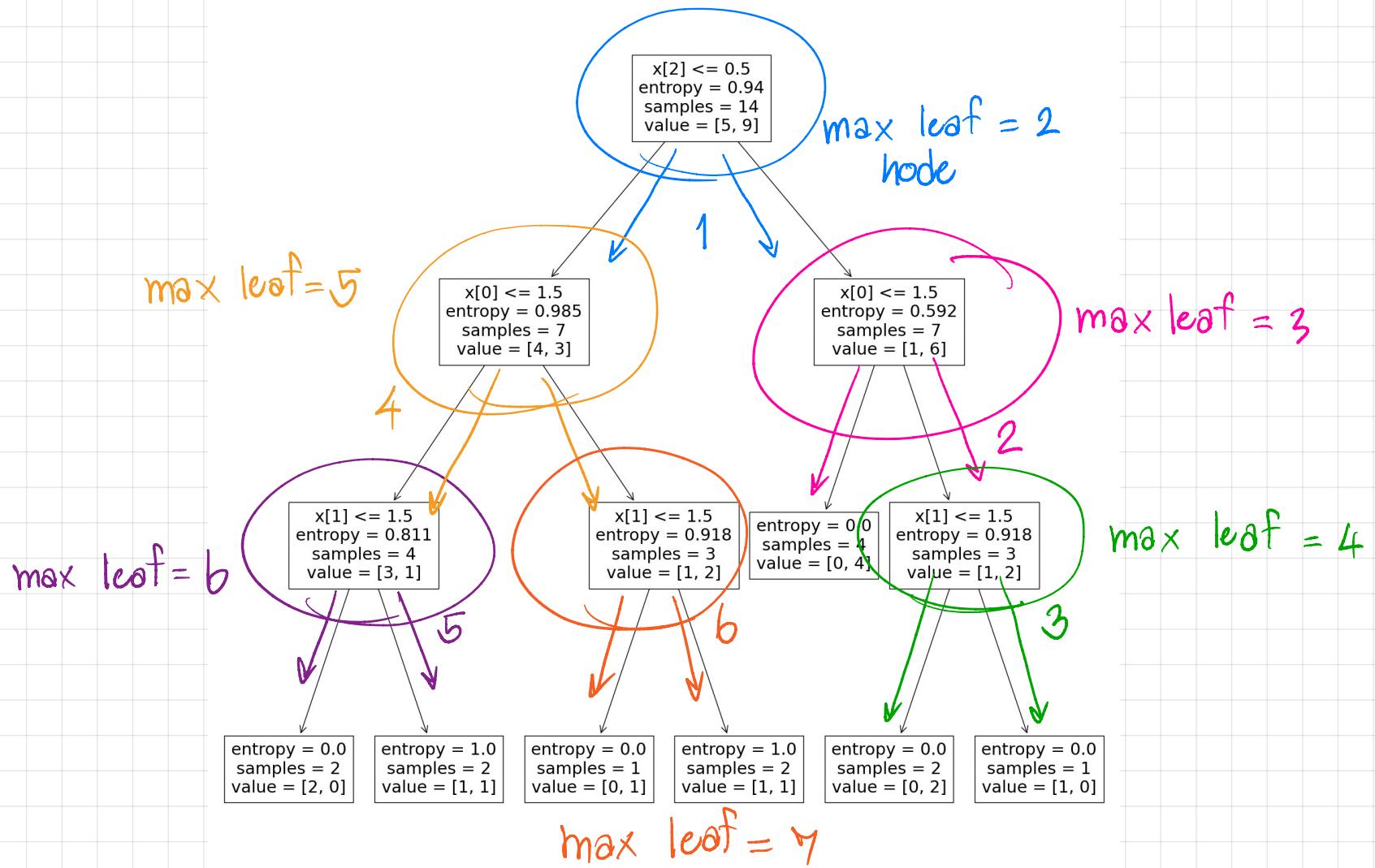
$$\text{Info}_{\text{student: no}}(D) = I(3, 4) = -\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right)$$
$$= 0.985$$

$$\text{Info}_{\text{age}}(D) = \frac{4}{7}I(1, 3) + \frac{3}{7}I(2, 1)$$
$$= \frac{4}{7}\left[-\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right)\right] + \frac{3}{7}\left[-\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right)\right]$$
$$= \frac{4}{7}(0.811) + \frac{3}{7}(0.918) = 0.857$$

$$\text{Info}_{\text{income}}(D) = \frac{3}{7}I(1, 2) + \frac{4}{7}I(2, 2)$$
$$= \frac{3}{7}\left[-\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right)\right] + \frac{4}{7}\left[-\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right)\right]$$
$$= \frac{3}{7}(0.918) +$$







# Decision Tree Induction: Algorithm

## □ Basic algorithm

ទូទាត់អនុវត្ត

បង្ហាញការពារេលយោ

□ Tree is constructed in a top-down, recursive, divide-and-conquer manner

feature

□ At start, all the training examples are at the root

ម៉ាកក់

□ Examples are partitioned recursively based on selected attribute

entropy ស្ថុ

□ 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

លើកសរសៃណីមែនបាន

ស្ថុលេខា  
entropy អាចកើតឡើង

□ All samples for a given node belong to the same class → pure data

□ There are no remaining attribute for further partitioning

□ There are no samples left

Node មួយឯ pure (ការបង្ការេះបែប)

□ Prediction

ការបង្ការេះបែប

50 : 50

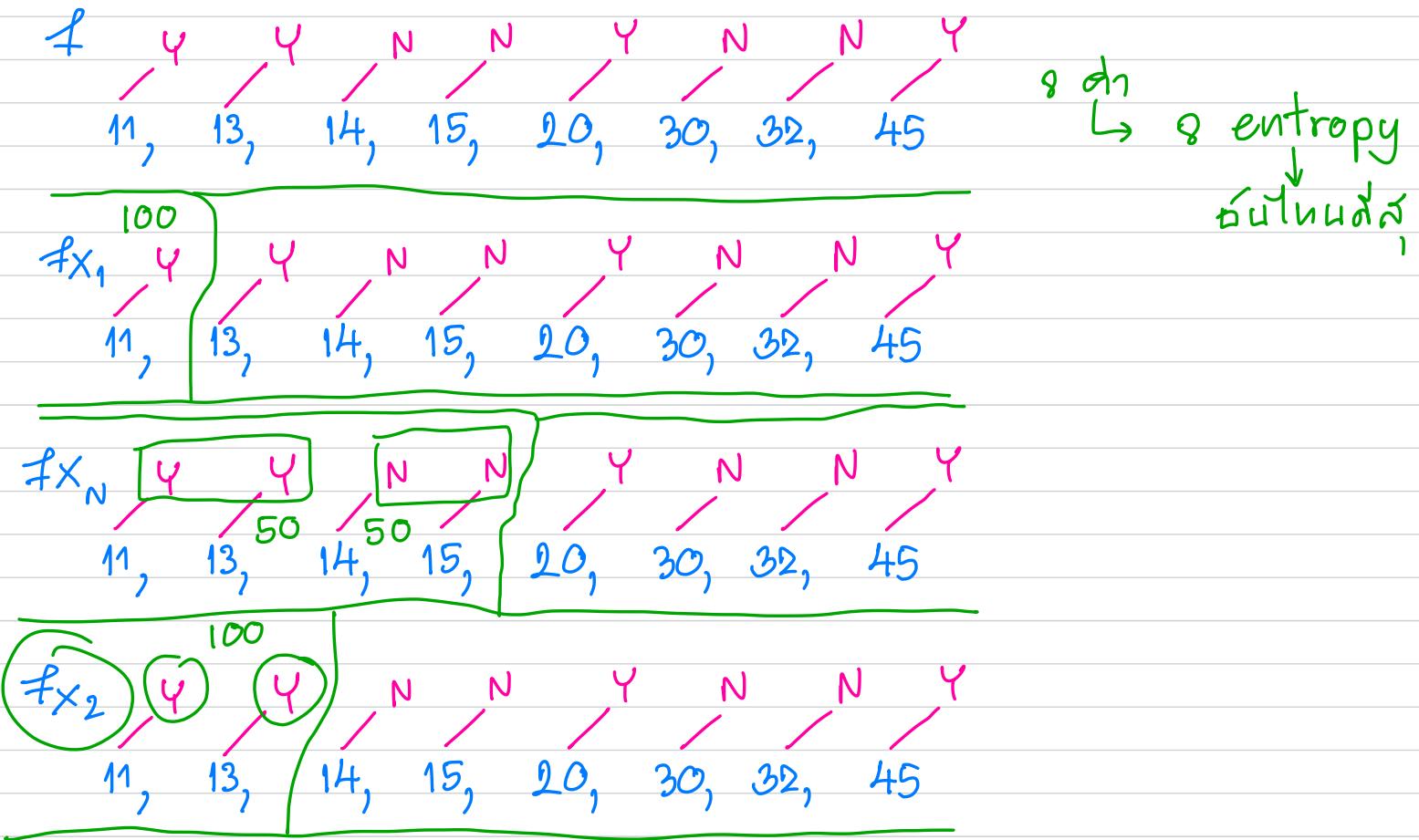
Y : N

□ **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  $A > P$



$Y = \text{pure}$   
ผิวนะ=เด่นๆ

กานต์  
→ ก้าเก็ตต้าเวย์  
continuous

↑ ดี  
↑ entropy  
↑ ชั่นไบเนติกส์,

# 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)

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

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

---

- Example: 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  $D_2$

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

- Gini<sub>{low,high}</sub> is 0.458; Gini<sub>{medium,high}</sub> 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 Three 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

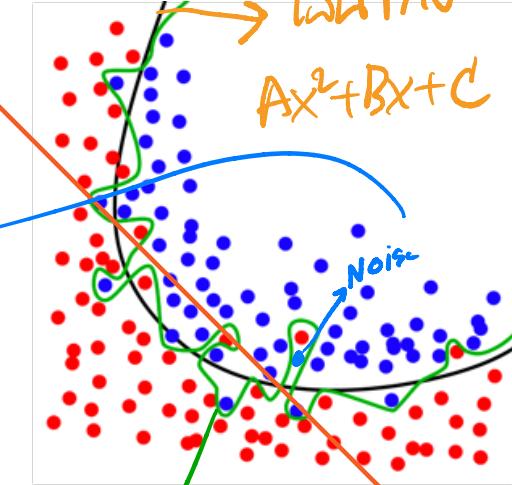
# Other Attribute Selection Measures

---

- Minimal Description Length (MDL) principle
  - Philosophy: The simplest solution is preferred
  - The best tree as the one that requires the fewest # of bits to both (1) encode the tree, and (2) encode the exceptions to the tree
- CHAID: a popular decision tree algorithm, measure based on  $\chi^2$  test for independence
- Multivariate splits (partition based on multiple variable combinations)
  - CART: finds multivariate splits based on a linear combination of attributes
- There are many other measures proposed in research and applications
  - E.g., G-statistics, C-SEP
- Which attribute selection measure is the best?
  - Most give good results, none is significantly superior than others

# Overfitting and Tree Pruning

- ❑ Overfitting: An induced tree may overfit the training data
- ❑ Too many branches, some may reflect anomalies due to noise or outliers
  - ຖົນທຳກັນ
- ❑ 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
  - ເປັນດູກອນ
  - ສູນສຸດ
  - $y = mx + c$
- ❑ Use a set of data different from the training data to decide which is the “best pruned tree”



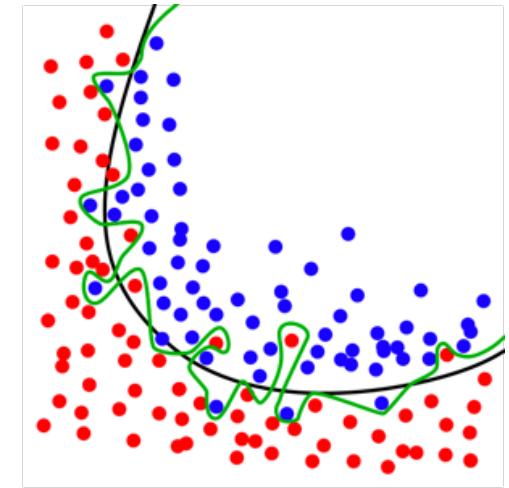
ກົດລວມກົມ  
ກົມໄດ້ກົມ

ກົມໄດ້ກົມ

ເປັນດູກອນ  
ສູນສຸດ  
 $y = mx + c$

# Overfitting and Tree Pruning

- 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”



## CODE

PACKAGE : Sklearn Classification

1. Input
  2. Define Model
  3. Training (.fit ( $X, Y$ ))
  4. test - ချက်ဆုပ် (.predict ( $\vec{X}$ ))  $\rightarrow \hat{Y}$
- Condition အေမြေပုဂ္ဂန်းနှင့် သတ်မှတ်မှု  
ထုတေသန  
ကိုယ် / လုပ်စဉ်

# Classification in Large Databases

---

- Classification—a classical problem extensively studied by statisticians and machine learning researchers
- Scalability: Classifying data sets with millions of examples and hundreds of attributes with reasonable speed
- Why is decision tree induction popular?
  - Relatively fast learning speed
  - Convertible to simple and easy to understand classification rules
  - Easy to be adapted to database system implementations (e.g., using SQL)
  - Comparable classification accuracy with other methods
- **RainForest** (VLDB'98 — Gehrke, Ramakrishnan & Ganti)
  - Builds an AVC-list (attribute, value, class label)

# RainForest: A Scalable Classification Framework

- ❑ The criteria that determine the quality of the tree can be computed separately
  - ❑ Builds an AVC-list: **AVC (Attribute, Value, Class\_label)**
- ❑ **AVC-set** (of an attribute  $X$ )
  - ❑ Projection of training dataset onto the attribute  $X$  and class label where counts of individual class label are aggregated
- ❑ **AVC-group** (of a node  $n$ )
  - ❑ Set of AVC-sets of all predictor attributes at the node  $n$

age	income	student	credit_rating	comp
<=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

The Training Data

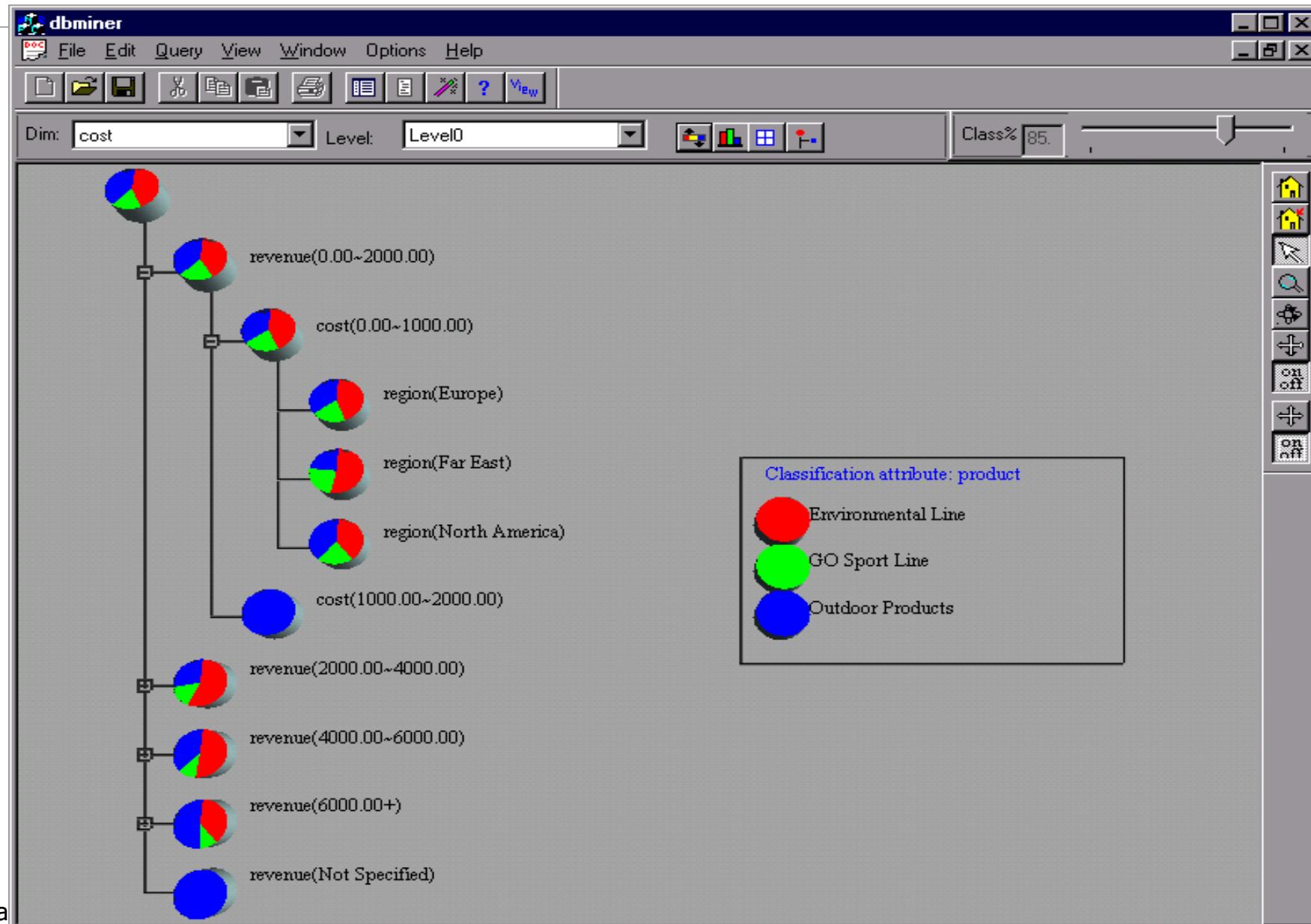
AVC-set on Age	
Age	Buy_Computer
yes	no
<=30	2
31..40	4
>40	3
	2

AVC-set on Income	
income	Buy_Computer
yes	no
high	2
medium	2
low	1

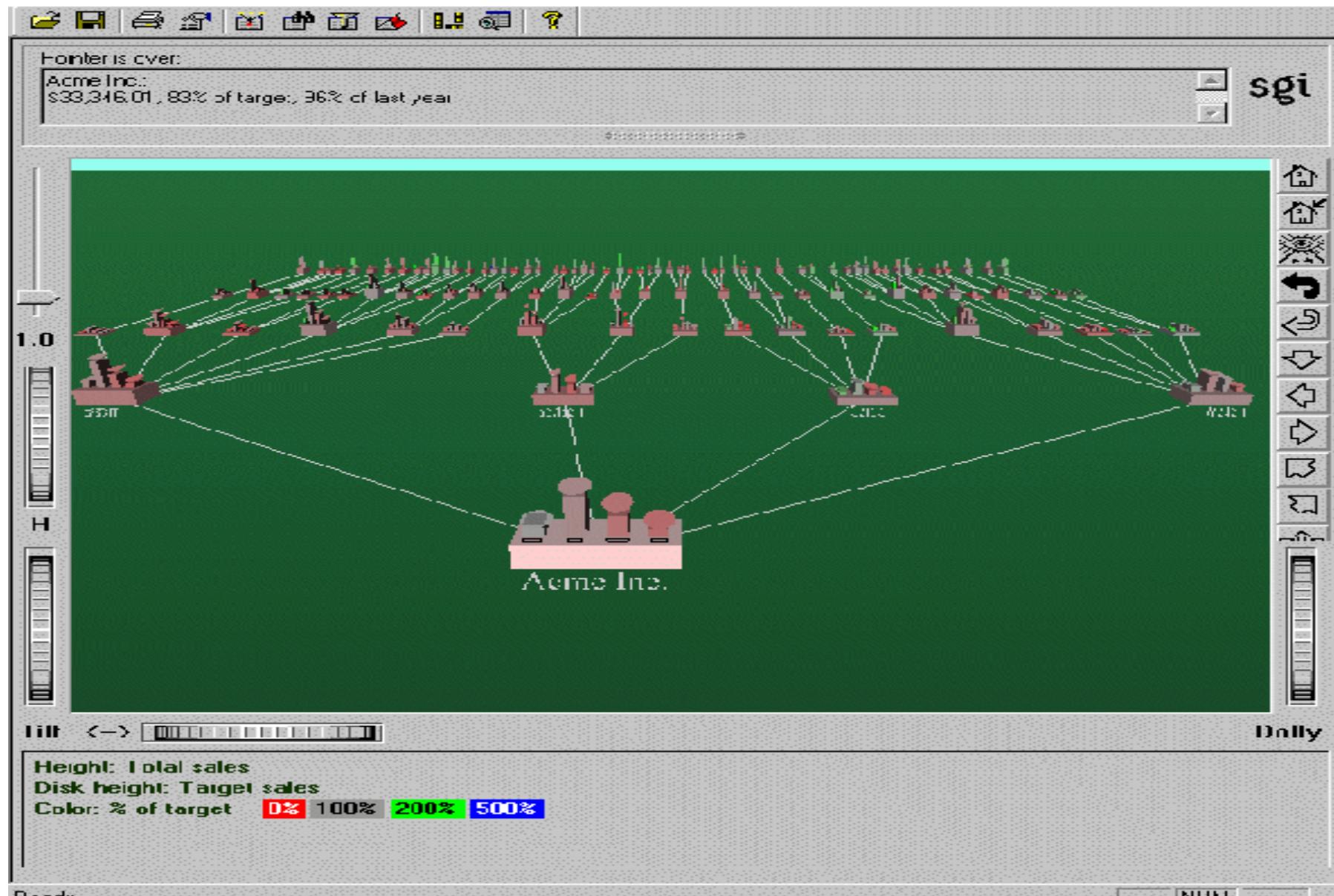
AVC-set on Student		AVC-set on Credit_Rating	
student	Buy_Computer	Credit rating	Buy_Computer
	yes		no
yes	6	fair	2
	1		6
no	3	excellent	3
	4		3

Its AVC Sets

# Presentation of Classification Results

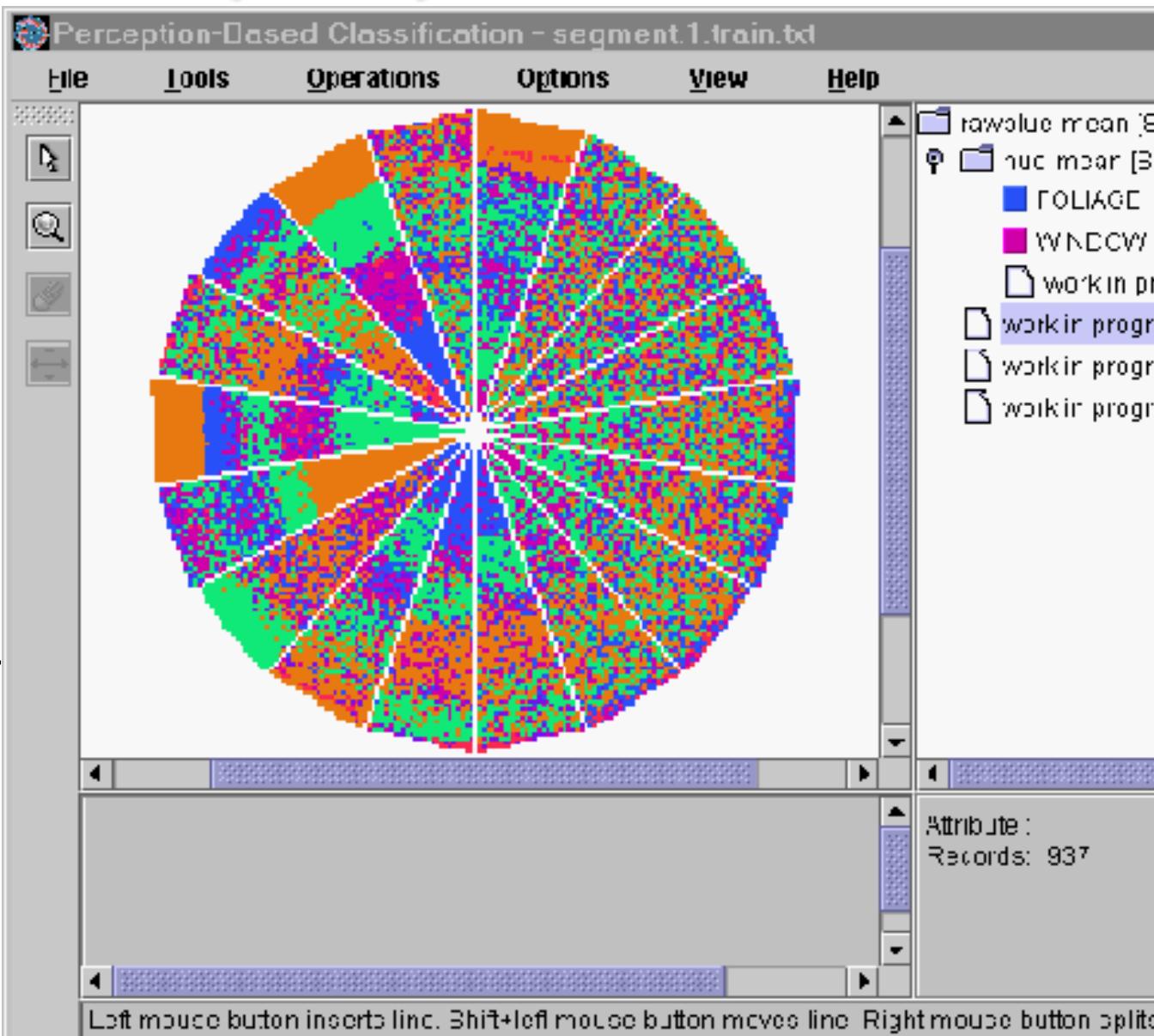


# Visualization of a Decision Tree (in SGI/MineSet 3.0)

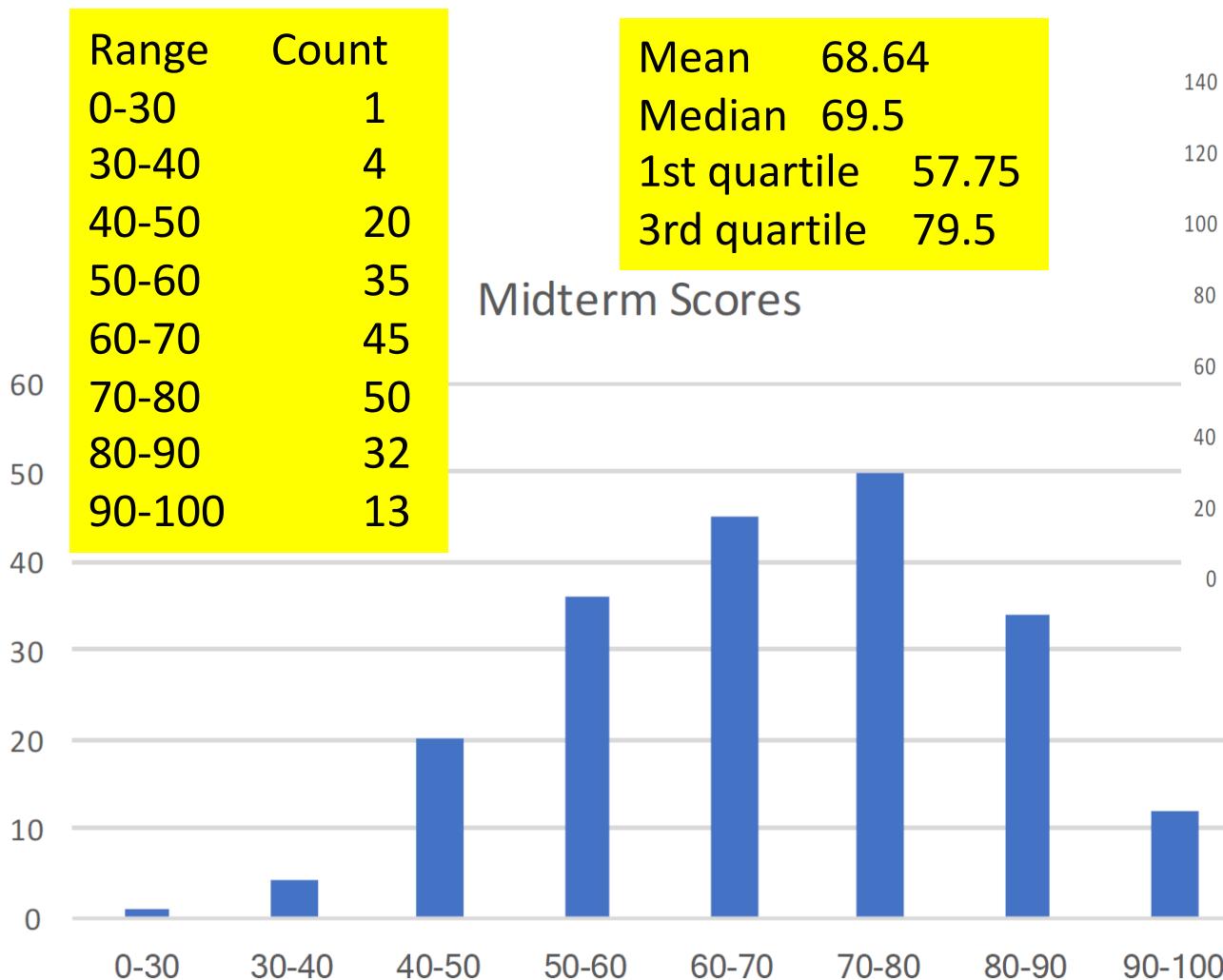


# Interactive Visual Mining by Perception-Based Classification (PBC)

- ❑ Perception-based classifier (PCB): developed at Univ. of Munich (1999)
- ❑ One color represents one class label
- ❑ One pie represents one attribute (or variable)
- ❑ The pie with random spread implies weak classification power
- ❑ The pie with clearly partitioned color strips implies good classification power
- ❑ One can select a good attribute and regenerate new pie charts for classification at the subsequent levels



# CS412-Fall 2017: Midterm Statistics



Midterm Options

