

# Supervised vs. Unsupervised Learning (1)

အသုံးပြုမှုအမျိုးအမည်  
သတ်မှတ်ထားသော

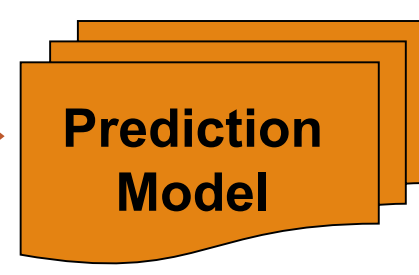
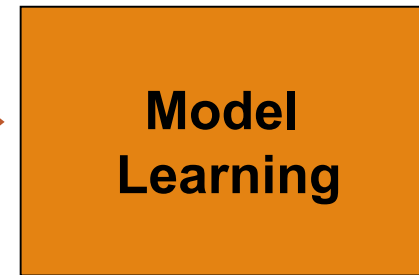
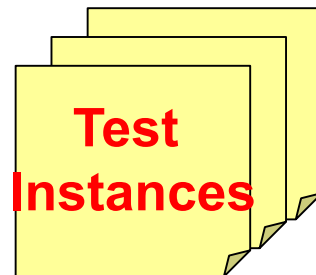
အသုံးပြုမှုအမျိုးအမည်  
သတ်မှတ်ထားသော

## 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

Training Data with class label:

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

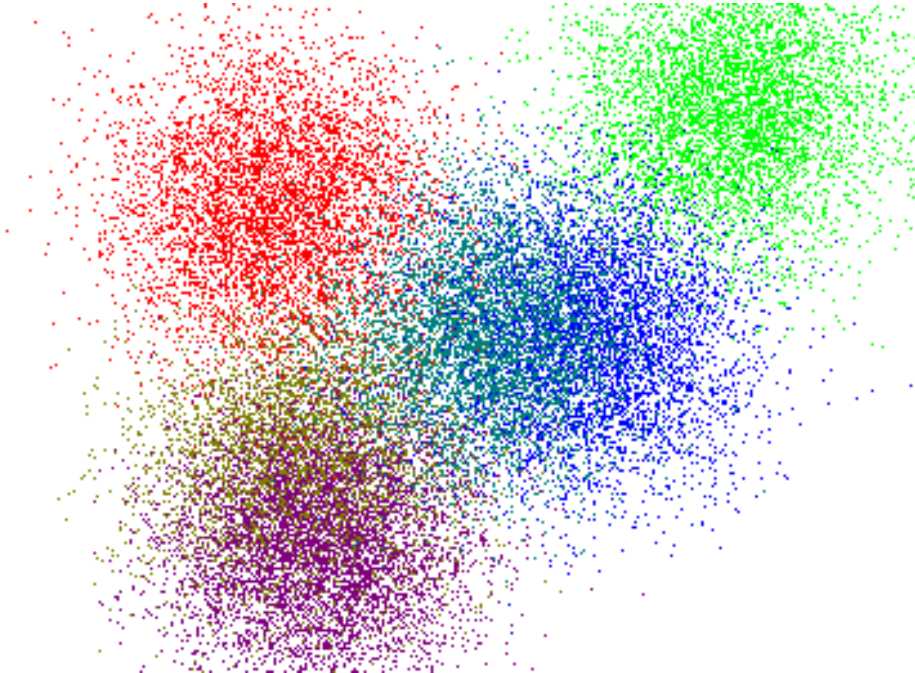


# Supervised vs. Unsupervised Learning (2)

## □ Unsupervised learning (clustering)

အဲဒါနဲ့ပတ်သက်တဲ့မေးခွန်းများကိုလည်းကောင်း

- 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

การหาค่าของฟังก์ชัน

การหาค่า Regression

## Classification

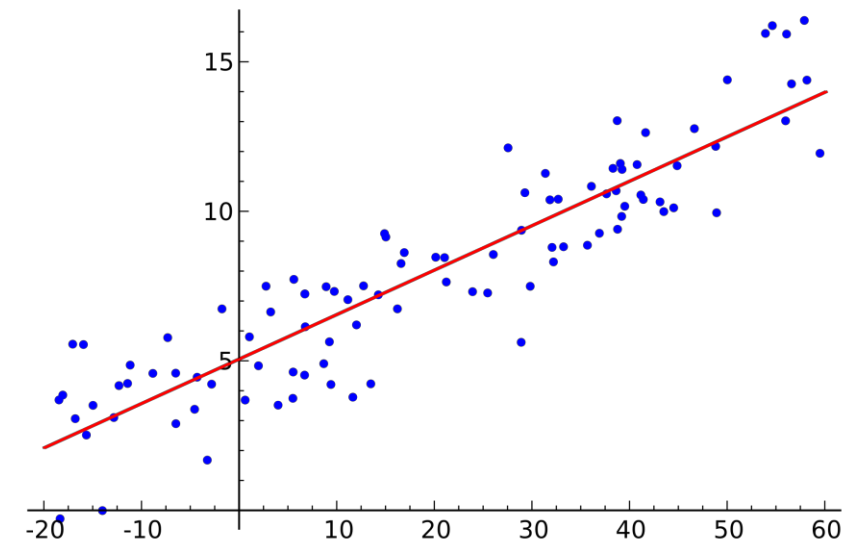
- Predict categorical class labels (discrete or nominal)
- Construct a model based on the training set and the **class labels** (the values in a classifying attribute) and use it in classifying new data

## 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



# Classification—Model Construction, Validation and Testing

- ❑ **Model construction** *เอา data ที่สลับที่กันแล้วเอาค่าของ data ที่มันอยู่ในโมเดลมาสร้างโมเดลโดยที่มันสลับที่*
  - ❑ 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:** *เอาโมเดลไปทดสอบ*
  - ❑ **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

# 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  $D$ :

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

- ❑ 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  $A$

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

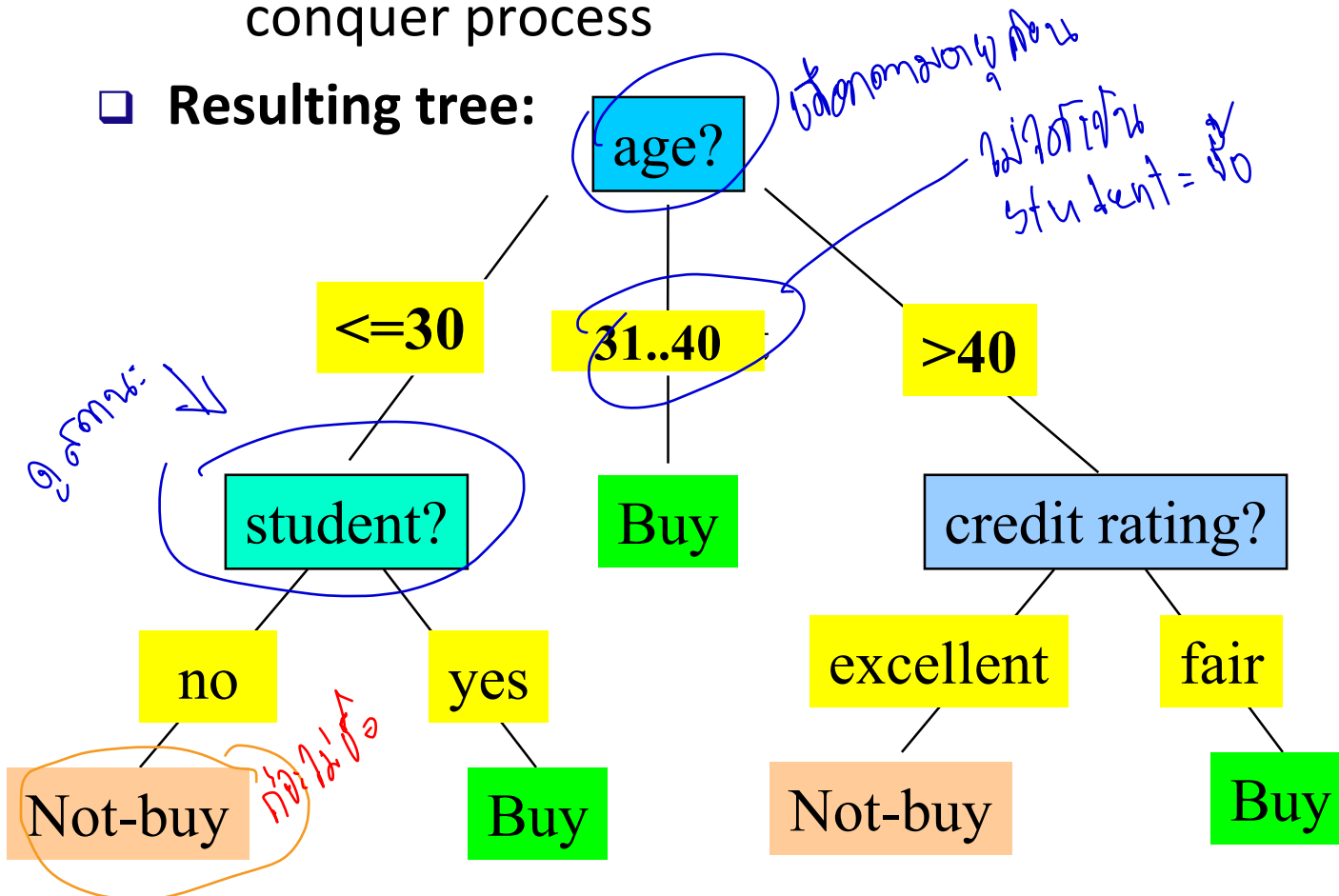


# Decision Tree Induction: An Example

## Decision tree construction:

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

## Resulting tree:



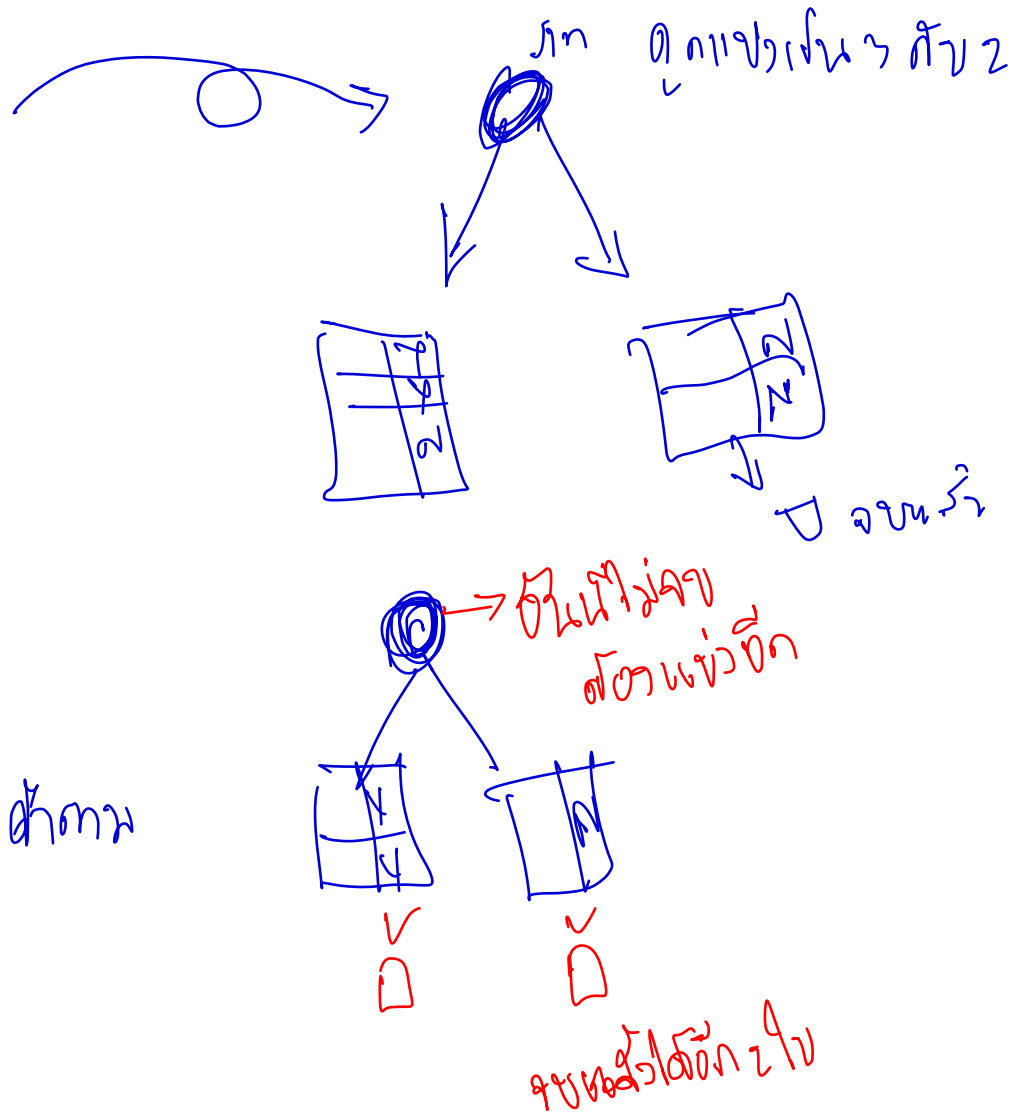
Training data set: Who buys computer?

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	<del>no</del>
<=30	high	no	excellent	<del>no</del>
31...40	high	no	fair	<del>yes</del>
>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

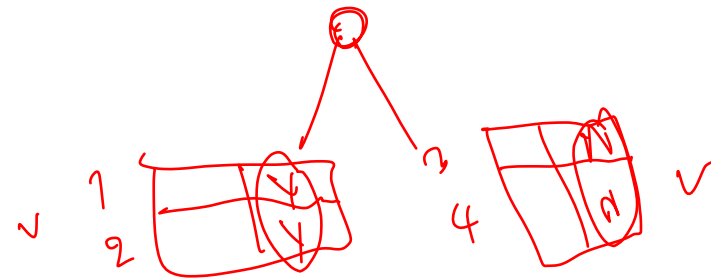
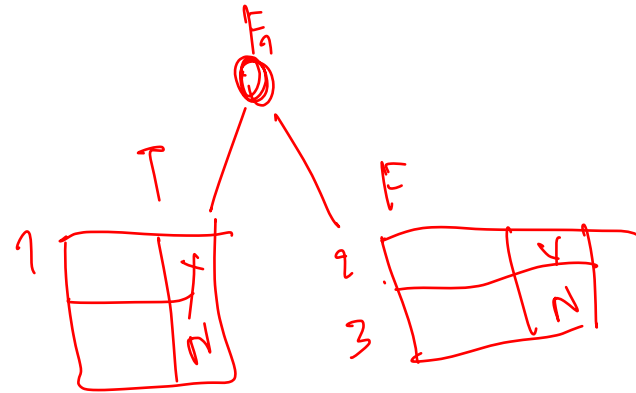
Note: The data set is adapted from "Playing Tennis" example of R. Quinlan

				Y
				Z
				Y
				Z
				Z

- สถานะของสีที่ปรากฏ
- จุดที่เปลี่ยนสีใน 3 สถานะ



$F_1$	$F_2$	$F_3$	$Y$
T	T	F	Y
F	T	F	Y
F	F	F	N
T	F	F	N



more True n yes

more false n no



# 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 <sub>i</sub>	n <sub>i</sub>	I(p <sub>i</sub> , n <sub>i</sub> )
<=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

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

*Handwritten notes: A stick figure points to the 'age' attribute. Above the first term, '≤30' is written with 'Y' and 'N' below it. Above the second term, '31-40' is written with 'Y' and 'N' below it. Below the third term, '740' is written.*

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

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

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