# CS 412 Intro. to Data Mining

## Chapter 8. Classification: Basic Concepts

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

1

# 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
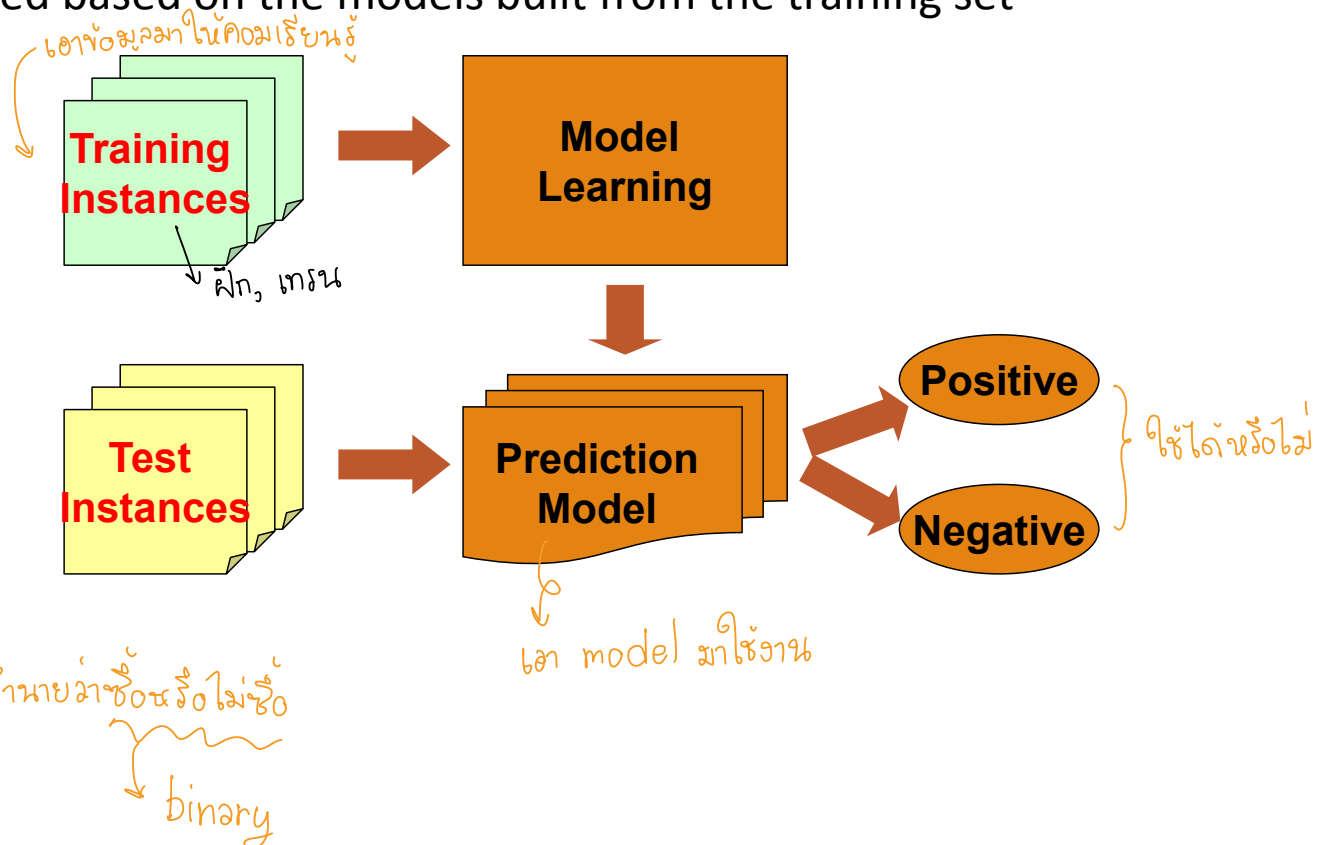
□ Additional Concepts on Classification

□ Summary

3

# Supervised vs. Unsupervised Learning (1)

❑ Supervised learning (classification)

มีผู้สอน

❑ 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
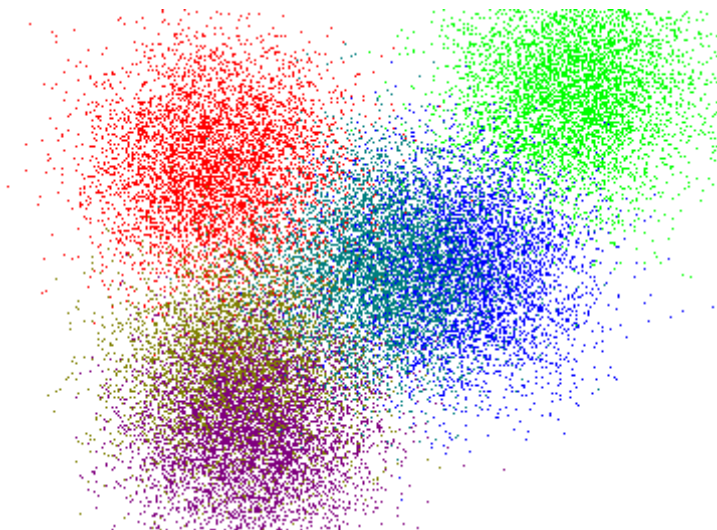


Training Data with class label:

| 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 |

เอาข้อมูลมาให้คอมเรียนรู้

Training Instances

ฝึก, เทรน

Model Learning

Test Instances

Prediction Model

เอา model มาใช้งาน

Positive

Negative

ใช้ได้หรือไม่

มีคำตอบ, ทำนายว่าซื้อหรือไม่ซื้อ

binary

4

# Supervised vs. Unsupervised Learning (2)

ไม่มีผู้สอน , ไม่มีจุดมุ่งหมายในการเรียน

❑ Unsupervised learning (clustering)   ไม่มีคำตอบ

  ❑ 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
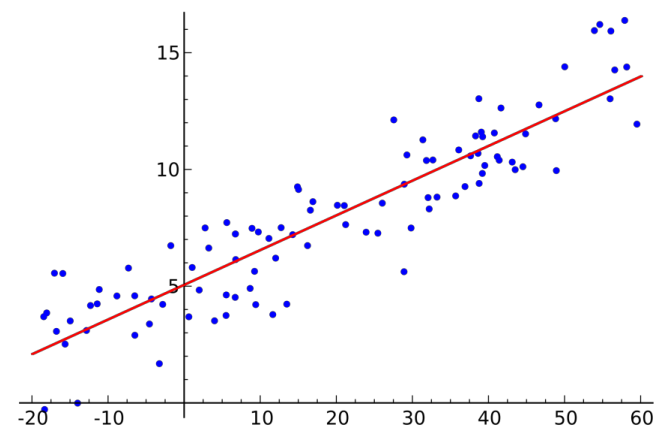
ใส่ข้อมูลแล้วให้มันแบ่งกลุ่ม

5

# Prediction Problems: Classification vs. Numeric Prediction

❑ Classification   → ทำนาย class, กลุ่มคน

   ❑ Predict categorical class labels (discrete or nominal)   ทำนายว่าอยู่กลุ่มไหน

   ❑ 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**
  - ❑ 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

8

ต้นไม้ตัดสินใจ

# Decision Tree Induction: An Example

ทำนายว่าใคร: ซื้อไม่ ซื้อ

X (Feature)    y (label)

- **Decision tree construction**:
  - A top-down, recursive, divide-and-conquer process
- **Resulting tree**:

ราก

age?

<=30    31..40    >40

student?    Buy    credit rating?

no    yes    excellent    fair

Not-buy    Buy    Not-buy    Buy

ใบ

9

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 |

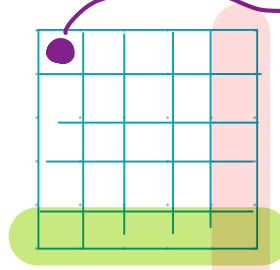Note: The data set is adapted from "Playing Tennis" example of R. Quinlan

Data Structure

$y = f(x)$

x ไปใส่ใน function แล้วเกิด y

# หลักการสร้าง Disition tree



x
y

1. ต้องเริ่มสร้างจากราก → กิ่ง → ... → ใบ

2. อย่าทำต้นไม้ให้สูง ควรทำให้เตี้ยที่สุด

ไม่ได้ใช้หลักการเงื่อนไขประพจน์   label จะไหลไปตามการตอบ

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

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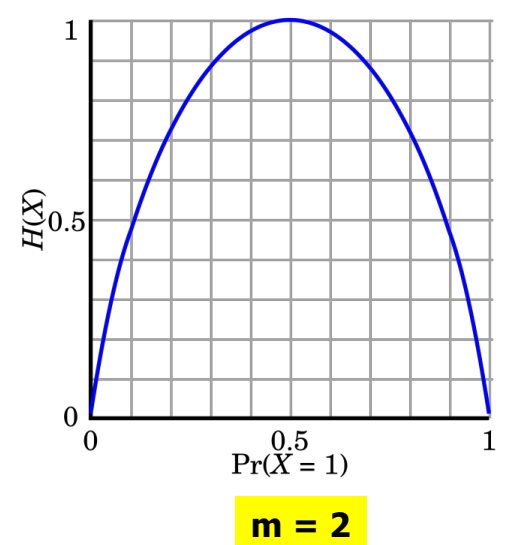  - ❑ Calculation: For a discrete random variable Y taking m distinct values $\{y_1, y_2, \ldots, y_m\}$

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

  - ❑ Interpretation
    - ❑ Higher entropy → higher uncertainty
    - ❑ Lower entropy → lower uncertainty
- ❑ Conditional entropy

$$H(Y|X) = \sum_{x} p(x)H(Y|X = x)$$

m = 2

# Information Gain: An Attribute Selection Measure

*เกณฑ์ทักของข้อมูล*

- ❑ 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:

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

- ❑ Information needed (after using A to split D into v partitions) to classify D:

$$Info_A(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times Info(D_j)$$    *คำนวณตามจำนวน Freature*

- ❑ Information gained by branching on attribute A

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

11

# Example: Attribute Selection with Information Gain

❑ Class P: buys_computer = "yes"

❑ Class N: buys_computer = "no"

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

*(yes, no)*

| 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 |

*14 row*

*h*

*S*

12

*<=30*  *31 - 40*

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

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

*>41*

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

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

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

$$I(A,B,C) = -\frac{A}{S}\log\frac{A}{S} - \dots - \frac{C}{S}\log\frac{C}{S}$$

*กรณีที่ไม่ใช้ label แค่ yes, no*

*เลือกตัวที่มาก สุด เป็นราก*

# Homework

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

Similarly, we can get

$Gain(income) = 0.029$

$Gain(student) = 0.151$

$Gain(credit\_rating) = 0.048$

① ②

เลือกตัวที่มาก สุด
เป็นราก

① พิจารณาเลือกราก
จากการคำนวณ

ดังนั้นๆได้

ซึ่งมี 3 กลุ่ม

age

$< 30$    $31 - 40$    $> 40$

มี 5     มี 4     มี 5

yes

ต่อมาคือ

Student

yes    No

# คำนวณ

$< 30$ ;

$Info(D) = I(2,3) = -\frac{2}{5}\log_2 \frac{2}{5} - \frac{3}{5}\log_2 \frac{3}{5}$  (yes, No)

$Info_{income}(D) = \frac{2}{5} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0)$  (high, medium, low)

$Info_{student}(D) = \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3)$  (yes, No)

$Info_{credit}(D) = \frac{3}{5} I(1,2) + \frac{2}{5}(1,1)$  (fair, exellent)

$31 \to 40$ ;

$Info(D) = I(4,0) = -\frac{4}{4}\log_2 \frac{4}{4} - \frac{0}{4}\log_2 \frac{0}{4}$  (yes, No)

$Info_{income}(D) = \frac{2}{4} I(1,1) + \frac{2}{5} I(0,1) + - I(1,0)$  (high, medium, low)

$Info_{student}(D) = \frac{2}{4} I(2,0) + \frac{2}{4} I(2,0)$  (yes, No)

$Info_{credit}(D) = \frac{2}{4} I(2,0) + \frac{2}{4}(2,0)$  (fair, exellent)

$> 40$ ;

$Info(D) = I(3,2) = -\frac{3}{5}\log_2 \frac{3}{5} - \frac{2}{5}\log_2 \frac{2}{5}$  (yes, No)

$Info_{income}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1)$  (medium, low)

$Info_{student}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1)$  (yes, No)

$Info_{credit}(D) = \frac{3}{5} I(3,0) + \frac{2}{5}(0,2)$  (fair, exellent)