




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) เป็นการสร้างโมเดลแบบมีผู้สอน

- 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

ข้อมูลลูกค้าร้านขายคอมพิวเตอร์

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

Training Instances

Model Learning

Test Instances

Prediction Model

มีค่าธรรมเนียม ซื้อหรือไม่ซื้อ เป็น
คุณสมบัติที่สนใจ ซึ่งแบบมีผู้สอนจะ
เป็นแบบมีคำตอบอยู่แล้ว ว่าซื้อ
หรือไม่ซื้อ เพื่อให้โปรแกรมเรียนรู้

Positive

Negative

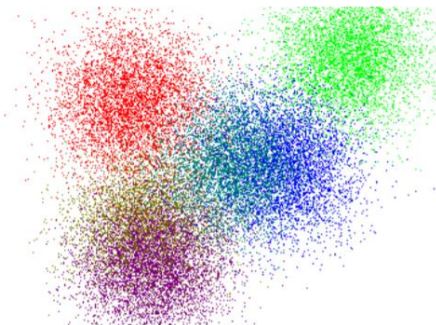
4

Supervised vs. Unsupervised Learning (2)

Unsupervised learning (clustering)

เป็นการสร้างโมเดลแบบไม่มีผู้สอน ไม่มี
จุดมุ่งหมายตั้งแต่ต้น เป็นเพียงการจัดกลุ่ม

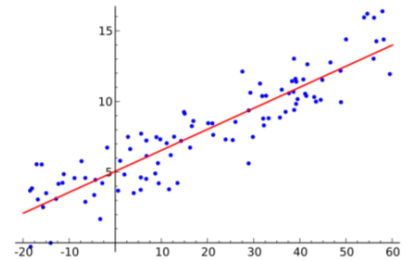
- The class labels of training data are unknown
- Given a set of observations or measurements, establish the possible existence of classes or clusters in the data



5

Prediction Problems: Classification vs. Numeric Prediction

- ❑ **Classification** สร้างโมเดลเพื่อใช้ทำนายค่าตอบ
 - ❑ 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



6

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

7

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

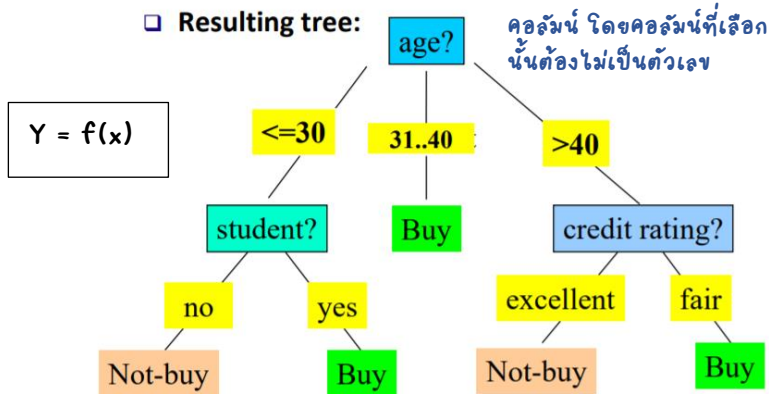
8

Decision Tree Induction: An Example

Decision tree construction:

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

Resulting tree:



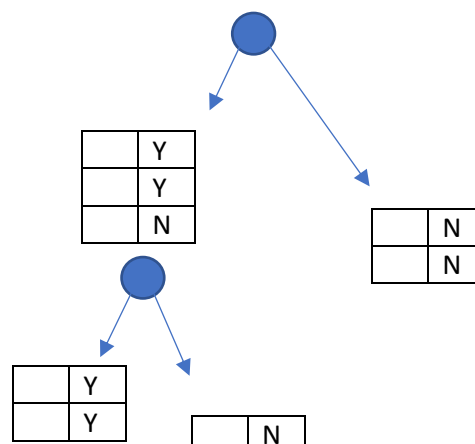
X = Feature				Y = Label
Training data set: Who buys computer?				
age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	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

9

				Y
				N
				Y
				N
				N

X 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 $\{y_1, y_2, \dots, y_m\}$

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

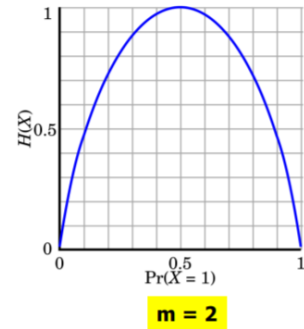
□ Interpretation

□ Higher entropy \rightarrow higher uncertainty

□ Lower entropy \rightarrow lower uncertainty

□ Conditional entropy

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



10

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

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

11

Example: Attribute Selection with Information Gain

□ Class P: buys_computer = "yes"

□ Class N: buys_computer = "no"

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

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'es and 3 no's.

Hence

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

Similarly, we can get

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$

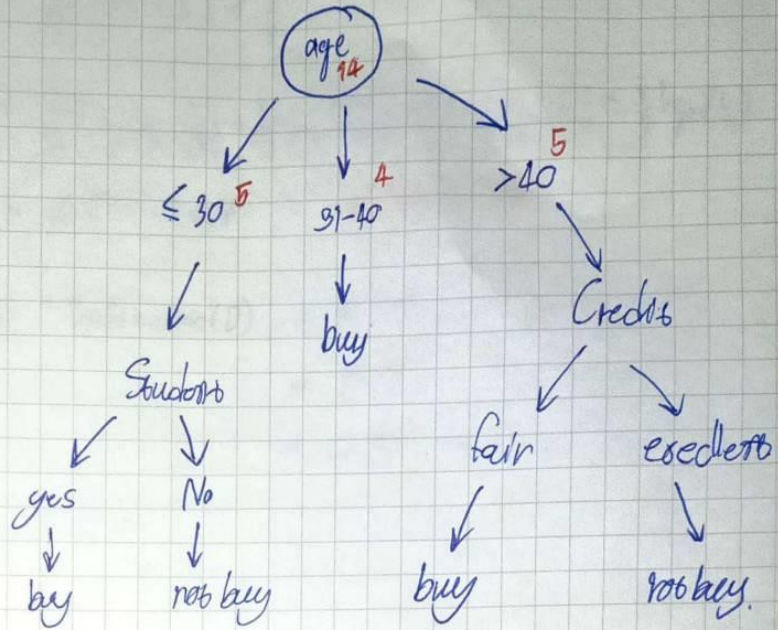
HW 14.

$$\text{Gain}(\text{age}) = 0.246$$

$$\text{Gain}(\text{Income}) = 0.029$$

$$\text{Gain}(\text{Students}) = 0.151$$

$$\text{Gain}(\text{Credit}) = 0.048$$



Age ≤ 30

$$\text{Info}(D) = I(2,3) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0.5288 + 0.4422 = 0.971$$

$$\begin{aligned} \text{Info income}(D) &= \frac{2}{5} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0) \\ &= \frac{2}{5} \left[-\frac{0}{2} \log_2\left(\frac{0}{2}\right) - \frac{2}{2} \log_2\left(\frac{2}{2}\right) \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right] \\ &\quad + \frac{1}{5} \left[-\frac{1}{1} \log_2\left(\frac{1}{1}\right) - \frac{0}{1} \log_2\left(\frac{0}{1}\right) \right] \\ &= 0 + \frac{2}{5} (0.5 + 0.5) + 0 = 0.40 \end{aligned}$$

$$\begin{aligned} \text{Info students}(D) &= \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3) \\ &= \frac{2}{5} \left[-\frac{2}{2} \log_2\left(\frac{2}{2}\right) - \frac{0}{2} \log_2\left(\frac{0}{2}\right) \right] + \frac{3}{5} \left[-\frac{0}{3} \log_2\left(\frac{0}{3}\right) - \frac{3}{3} \log_2\left(\frac{3}{3}\right) \right] \\ &= 0 - 0 = 0 \end{aligned}$$

Info credit (D)

$$= \frac{2}{5} I(1,2) + \frac{2}{5} I(1,2)$$

$$= \frac{2}{5} \left[-\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right]$$

$$= 0.5510 + 0.4$$

$$= 0.9510$$

Gain (Income) = Info (D) - Info income (D) = 0.951 - 0.40 = 0.551

Gain (Students) = 0.951 - 0 = 0.951

Gain (Credits) = 0.951 - 0.951 = 0.00

Age 31-40

Info (D) = I(4,0)

$$= -\frac{4}{4} \log_2 \left(\frac{4}{4} \right) - \frac{0}{4} \log_2 \left(\frac{0}{4} \right)$$

$$= 0$$

Info income (D)

$$= \frac{2}{4} I(2,0) + \frac{1}{4} I(1,0) + \frac{1}{4} I(1,0)$$

$$= \frac{2}{4} \left[-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{0}{2} \log_2 \left(\frac{0}{2} \right) \right] + \frac{1}{4} \left[-\frac{1}{1} \log_2 \left(\frac{1}{1} \right) - \frac{0}{1} \log_2 \left(\frac{0}{1} \right) \right]$$

$$+ \frac{1}{4} \left[-\frac{1}{1} \log_2 \left(\frac{1}{1} \right) - \frac{0}{1} \log_2 \left(\frac{0}{1} \right) \right] = 0$$

Info students (D)

$$= \frac{2}{4} I(2,0) + \frac{2}{4} I(2,0)$$

$$= \frac{2}{4} \left[-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{0}{2} \log_2 \left(\frac{0}{2} \right) \right] + \frac{2}{4} \left[-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{0}{2} \log_2 \left(\frac{0}{2} \right) \right]$$

$$= 0$$

Info credits (D)

$$= \frac{2}{4} I(2,0) + \frac{2}{4} I(2,0)$$

$$= \frac{2}{4} \left[-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{0}{2} \log_2 \left(\frac{0}{2} \right) \right] + \frac{2}{4} \left[-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{0}{2} \log_2 \left(\frac{0}{2} \right) \right]$$

$$= 0$$

Age 31-40 buy-computer = yes missing

Age > 40

$$\begin{aligned} \text{Info}(D) &= I(3,2) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \\ &= 0.4422 + 0.5288 = 0.9710 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{income}}(D) &= \frac{0}{5} I(0,0) + \frac{2}{5} I(2,1) + \frac{2}{5} I(1,1) \\ &= \frac{2}{5} \left[-\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right] \\ &= 0.551 + 0.4 = 0.951 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{Student}}(D) &= \frac{2}{5} I(2,1) + \frac{2}{5} I(1,1) \\ &= \frac{2}{5} \left[-\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right] \\ &= 0.551 + 0.4 = 0.951 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{Credit}}(D) &= \frac{2}{5} I(3,0) + \frac{2}{5} I(0,2) \\ &= \frac{2}{5} \left[-\frac{3}{3} \log_2\left(\frac{3}{3}\right) - \frac{0}{3} \log_2\left(\frac{0}{3}\right) \right] + \frac{2}{5} \left[-\frac{0}{2} \log_2\left(\frac{0}{2}\right) - \frac{2}{2} \log_2\left(\frac{2}{2}\right) \right] \\ &= 0 \end{aligned}$$

$$\text{Gain}(\text{Income}) = 0.020$$

$$\text{Gain}(\text{Student}) = 0.020$$

$$\text{Gain}(\text{Credit}) = 0.971$$