



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)

สร้าง model หรือ ใช้โมเดลทำนายค่าตอบ

- Supervised learning (classification)

- Supervision: The training data such as ^{มีคุณสมบัติ, มีค่าของหลายๆ} observations or ^{มีการสังเกต} measurements are accompanied by **labels** indicating the classes which they belong to
- ^{สร้างเพื่อเก็บข้อมูลจากหลายๆ column ของ file} 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

Training
Instances

Model
Learning

Data แบ่งออกเป็น
x ที่เราใช้
y (คำตอบ)

Test
Instances

Prediction
Model

Positive

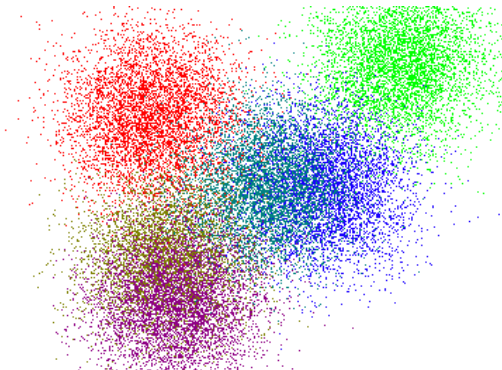
Negative

Supervised vs. Unsupervised Learning (2) ฝึกได้ x

- Unsupervised learning (clustering) ไม่ต้องใช้ข้อมูล, ไม่ต้องกำหนดคำตอบ

ฝึกอะไร → เสร็จไปไหน → ฝึกต่อกัน
ไปเรื่อย ๆ → ไปไหน

- 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

ทำนาย Class ไม่ถูก/ใช่

- **Classification**

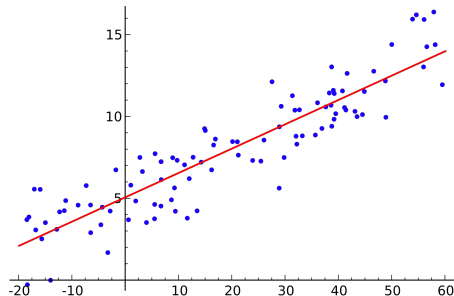
- Predict categorical class labels (discrete or nominal)
- Construct a model based on the training set and the **class labels** (the values in a classifying attribute) and use it in classifying new data

- **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

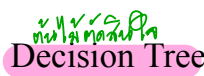


Classification—Model Construction, Validation and Testing

สร้าง model → วัดผล, ทำความเข้าใจ → นำไปใช้งาน

- **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

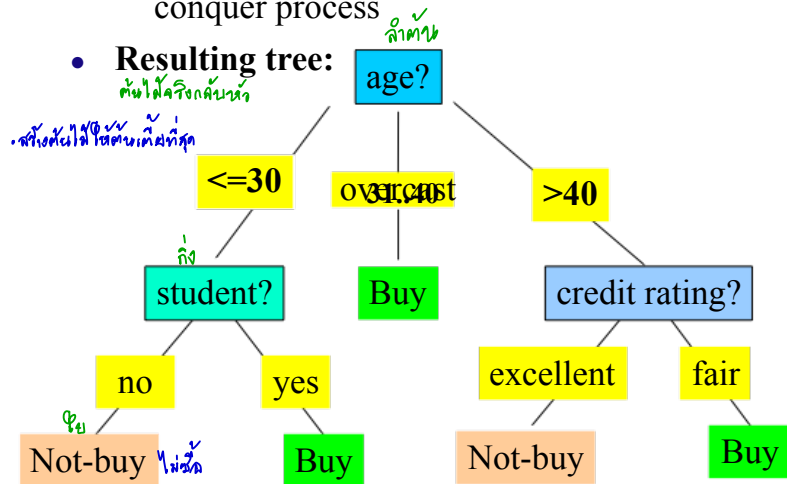
- Classification: Basic Concepts
-  Decision Tree Induction
- Bayes Classification Methods
- Linear Classifier
- Model Evaluation and Selection
- Techniques to Improve Classification Accuracy: Ensemble Methods
- Additional Concepts on Classification
- Summary

Decision Tree Induction: An Example

- Decision tree construction:**

- 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	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

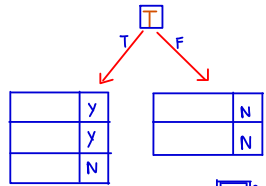
Note: The data set is adapted from
“Playing Tennis” example of R. Quinlan

ผล
รวม

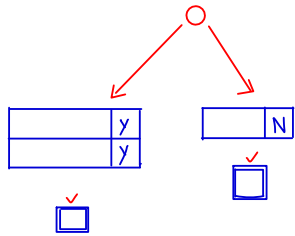
- สืบจาก root ก่อนคือขวาใน
- มี Data 2 ส่วน $\rightarrow x, y$
- เอา Data 5 ตัวมาแบ่ง root node (ตัวที่แบ่งได้ที่ดีที่สุด)

5 rows, 5 columns

				y
				N
				y
				N
				N



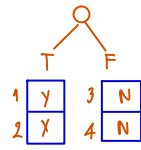
ใบแบ่งได้



เพิ่มอีก

	f_1	f_2	f_3	y
1	T	T	F	y
2	F	T	F	y
3	F	F	F	N
4	T	F	T	N

f_2



ไม่ผ่านใบต่อ

ดู f_1



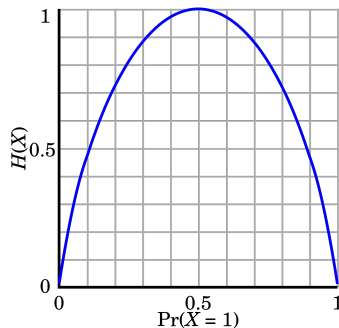
$T \rightarrow y$ 19 ตัว
 $F \rightarrow y$ 19 ตัว

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 $\{y_1, y_2, \dots,$

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

- Interpretation
 - Higher entropy \rightarrow higher uncertainty
 - Lower entropy \rightarrow lower uncertainty
- $C H(Y|X) = \sum_x p(x) H(Y|X = x)$



m = 2

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 :

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

pop
แต่ละชิ้นของต้นไม้ Info(D) จะน้อยลงถ้า

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

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

คำนวณตามต้นไม้ที่สร้าง

- Information gained by branching on attribute A

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

$$I(p_1, p_2, c) = -\frac{p_1}{S} \log_2 \frac{p_1}{S} - \frac{p_2}{S} \log_2 \frac{p_2}{S} - \frac{c}{S} \log_2 \frac{c}{S} \quad (\text{คำนวณของทุกตัว})$$

Example: Attribute Selection with Information Gain

- Class P: buys_computer = “yes”
- Class N: buys_computer = “no”

yes, no

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

แบ่ง yes กับ no ที่นี้ก็คือ

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

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

≤30
31-40
7 4 1

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

จะเปลี่ยนไปก็จริง

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$

Homework

Gain(age)

Age		Gain	
<=30	31-40	>40	
N	Y	Y	
N	Y	Y	
N	Y	N	
Y	Y	Y	
Y		N	

$$\text{Info}(p) = I(4,5) = -\frac{4}{14} \log_2 \frac{4}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.940$$

Gain(age)

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

$$\text{Info}_{\text{age}}(p) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

$$\text{Gain}(\text{age}) = 0.940 - 0.694 = 0.246$$

Gain(income)

Income	p_i	n_i	$I(p_i, n_i)$
high	2	2	1
medium	4	2	0.918
low	3	1	0.811

$$\text{Info}_{\text{income}}(p) = \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1) = 0.911$$

$$\text{Gain}(\text{income}) = 0.940 - 0.911 = 0.029$$

Gain(Student)

Student	p_i	n_i	$I(p_i, n_i)$
yes	6	1	0.992
no	3	4	0.985

$$\text{Info}_{\text{student}}(p) = \frac{7}{14} I(6,1) + \frac{7}{14} I(3,4) = 0.769$$

$$\text{Gain}(\text{student}) = 0.940 - 0.769 = 0.171$$

Gain(credit_rating)

credit_rating	p_i	n_i	$I(p_i, n_i)$
fair	6	2	0.811
excellent	3	3	1

$$\text{Info}_{\text{credit_rating}}(p) = \frac{7}{14} I(6,2) + \frac{7}{14} I(3,3) = 0.892$$

$$\text{Gain}(\text{credit_rating}) = 0.940 - 0.892 = 0.048$$

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

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

$$\text{Gain}(\text{Student}) = 0.171$$

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

größtmögliche Gain(age) ist der root node

if age <= 30

$$\text{Info}(p) = I(2,3) = 0.971$$

$$\text{Info}_{\text{income}}(p) \text{ vs age } (<=30)$$

Income	p_i	n_i	$I(p_i, n_i)$
high	0	2	0
medium	1	1	1
low	1	0	0

$$\text{Info}_{\text{income}}(p) \text{ vs age } (<=30) = \frac{2}{5} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0) = 0.4$$

$$\text{Gain}(\text{income}) \text{ vs age } (<=30) = 0.971 - 0.4 = 0.571$$

Info_student(p)

$$\text{Info}_{\text{student}}(p) \text{ vs age } (<=30) = \frac{2}{5} I(2,0) + \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3)$$

if age <= 30
 yes → yes (buy_computer) → 1
 no → no (buy_computer) → 0

age (>40)

$$\text{Info}(D) = I(3,2) = 0.971$$

$$\text{Info}_{\text{income}}(D) \text{ vs age } (>40)$$

Income	p_i	n_i	$I(p_i, n_i)$
low	1	1	1
medium	2	1	0.918

$$\text{Info}_{\text{income}}(D) \text{ vs age } (>40) = \frac{1}{3} I(1,1) + \frac{2}{3} I(2,1)$$

$$= 0.951$$

$$\text{Gain}(\text{income}) \text{ vs age } (>40) = 0.971 - 0.951$$

$$= 0.020$$

Info_{student} vs age (>40)

Student	p_i	n_i	$I(p_i, n_i)$
yes	2	1	0.916
no	1	1	1

$$\text{Info}_{\text{student}}(D) \text{ vs age } (>40) = \frac{2}{3} I(2,1) + \frac{1}{3} I(1,1) = 0.951$$

$$\text{Gain}(\text{student}) \text{ vs age } (>40) = 0.971 - 0.951 = 0.020$$

$$\text{Info}_{\text{credit-rating}}(D) \text{ vs age } (>40)$$

$$\text{Info}_{\text{credit-rating}}(D) \text{ vs age } (>40) = \frac{3}{5} I(3,0) + \frac{2}{5} I(0,2)$$

fair \rightarrow yes $\frac{3}{5}$
 excellent \rightarrow No $\frac{2}{5}$

